

# A Mathematical Introduction to Causality

**Lecture Notes**

Dr. Patrick Forré

Joris M. Mooij

September 25th, 2023

# Foreword

Causality is a broad topic, and these lecture notes cover only part of it. They originated over a period of three years as a by-product of the course on causality we taught for MSc. mathematics students. Our aim was to give a mathematically rigorous exposition of the graphical account to causal modeling, reasoning and inference, in the spirit of Wright, Spirtes, Glymour, Scheines, Pearl, and many others. Since there seemed to be no book or lecture notes out there that would fit our purpose, we decided to write our own.

The amount of material has grown over the years, and is still growing. We treat causal modeling with causal Bayesian networks (also known as ‘DAGs’) and structural causal models. Some unique features of our exposition are:

1. we have extended the standard formalisms with *input nodes* to enable a measure-theoretically rigorous treatment of the families of probability distributions that result from perfect interventions;
2. we allow for (sufficiently weak) cycles in structural causal models.

Our treatment is self-contained: we start with the basic definitions (with as prerequisites only basic measure theory and probability theory), and derive everything that is necessary to prove the validity of Markov properties, the do-calculus, adjustment criteria, all the way up to extended versions of the ID algorithm and the FCI algorithm. We show how—with relatively little extra work—the framework of causal modeling with directed acyclic graphs can be extended to directed graphs that may have cycles.

While the advantages of mathematical rigor should be obvious, the price paid is that the non-trivial conceptual issues are sometimes clouded by technicalities. We believe that our treatment fills a much needed gap in the literature on causality, and consider it complementary to the many existing writings on similar topics (which often focus more on concepts and less on mathematical rigor).

We are indebted to our teaching assistants Leon Lang, Philip Boeken, Pim de Haan and Noud de Kroon for providing feedback and for spotting several errors in earlier drafts. While the current version undoubtedly still contains mistakes, we believe that it is now ready for wider exposure. We appreciate any feedback that the reader may have, be it on content, typos, or (we hope not) more serious mistakes.

Joris Mooij & Patrick Forré  
Amsterdam  
June 2023

# Table of Contents

<b>Foreword</b>	<b>2</b>
<b>Contents</b>	<b>6</b>
<b>1. Experimental Causal Discovery</b>	<b>7</b>
1.1. Types of Correlations . . . . .	7
1.2. Causal Effects in the Real World . . . . .	8
1.3. Randomized Controlled Trials (RCT) . . . . .	10
<b>2. Transition Probability Theory</b>	<b>12</b>
2.1. Elementary Probability Theory . . . . .	12
2.2. Recap - Measure Theoretic Probability . . . . .	14
2.2.1. Measurable Spaces and Maps . . . . .	14
2.2.2. Finite and Probability Measures . . . . .	15
2.2.3. The Measure Integral . . . . .	16
2.2.4. The Lebesgue Measure . . . . .	17
2.3. Transition Measures and Markov Kernels . . . . .	17
2.3.1. Core Definitions . . . . .	17
2.3.2. Special Cases of Markov Kernels . . . . .	20
2.3.3. The Doob-Radon-Nikodym Derivative . . . . .	21
Proofs - Theorem of Doob-Radon-Nikodym . . . . .	26
2.3.4. Transition Probability Spaces . . . . .	29
2.4. Constructing Markov Kernels from Others . . . . .	30
2.4.1. Marginal Markov Kernels . . . . .	30
2.4.2. Product of Markov Kernels . . . . .	30
2.4.3. Composition of Markov Kernels . . . . .	32
2.4.4. Push-Forward of Markov Kernels . . . . .	33
2.4.5. Conditional Markov Kernels . . . . .	34
Proofs - Disintegration of Markov Kernels . . . . .	36
2.5. Conditional Independence . . . . .	41
2.5.1. Independence for Random Variables . . . . .	41
2.5.2. Conditional Independence for Random Variables . . . . .	43
2.5.3. Conditional Independence for Conditional Random Variables . . . . .	45
2.5.4. Example: Linear Gaussian Markov Kernels . . . . .	55
2.6. Separoid Axioms for Conditional Independence . . . . .	57
Proofs - Separoid Axioms for Conditional Independence . . . . .	59
2.7. Markov Kernels from Deterministic Mappings . . . . .	65
Proofs - Deterministic Representation of Markov Kernels . . . . .	68
<b>3. Graph Theory</b>	<b>75</b>
3.1. Core Concepts . . . . .	75

3.2. Operations on Graphs . . . . .	79
3.2.1. Hard Interventions on Graphs . . . . .	79
3.2.2. Node Splitting Interventions on Graphs . . . . .	81
3.2.3. Intervention Nodes . . . . .	82
3.2.4. Marginalization of Graphs . . . . .	84
3.3. $\sigma$ -Separation . . . . .	85
Proofs - $\sigma$ -Open Walks and Paths . . . . .	88
3.4. $d$ -Separation . . . . .	90
3.5. Acyclifications . . . . .	91
3.6. Separoid Axioms for $\sigma$ -/ $d$ -Separation . . . . .	94
Proofs - Separoid Axioms for $\sigma$ -/ $d$ -Separation . . . . .	96
<b>4. Causal Bayesian Networks</b>	<b>100</b>
4.1. Core Concepts . . . . .	100
4.2. Global Markov Property . . . . .	103
Proofs - Global Markov Property . . . . .	104
4.3. Operations on Causal Bayesian Networks . . . . .	111
4.3.1. Hard Interventions on Causal Bayesian Networks . . . . .	111
4.3.2. Node-Splitting Hard Interventions on Causal Bayesian Networks . . . . .	112
4.3.3. Soft Interventions on Causal Bayesian Networks . . . . .	112
4.3.4. Marginalization of Causal Bayesian Networks . . . . .	113
4.4. Standard Forms of Causal Bayesian Networks . . . . .	114
Proofs - Standard Forms of Causal Bayesian Networks . . . . .	116
<b>5. Identification of Causal Effects in CBNs</b>	<b>118</b>
5.1. Do-Calculus . . . . .	118
Proofs - Do-Calculus . . . . .	123
5.2. Adjustment Criteria and Formulae . . . . .	132
5.3. The ID-Algorithm . . . . .	137
5.3.1. Core Definitions and Notations . . . . .	137
5.3.2. The Interventional Ordered Local Markov Property . . . . .	141
Proofs - The Interventional Ordered Local Markov Property . . . . .	142
5.3.3. Ancestral Sets and Districts . . . . .	145
5.3.4. The ID-Algorithm . . . . .	147
Examples . . . . .	150
Proofs - Soundness Criteria . . . . .	156
<b>6. Structural Causal Models</b>	<b>164</b>
6.1. Motivation . . . . .	164
6.2. Solving SCMs . . . . .	165
6.2.1. Potential Outcomes . . . . .	167
6.2.2. Solutions . . . . .	168
6.2.3. Markov Kernels of Solutions . . . . .	170
6.2.4. Solution functions . . . . .	170

6.3. Interventions . . . . .	172
6.3.1. Hard interventions . . . . .	173
6.3.2. Soft interventions (mechanism changes) . . . . .	175
6.3.3. Intervention variables . . . . .	176
6.4. Composition and decomposition . . . . .	176
6.5. Unique solvability and simple SCMs . . . . .	177
6.6. Equivalences . . . . .	181
6.7. Marginalizations . . . . .	184
6.8. Graphs of SCMs . . . . .	188
6.9. Interventional equivalence of acyclic SCMs and (L-)CBNs . . . . .	192
6.10. Examples . . . . .	194
<b>7. Markov property for simple SCMs</b>	<b>199</b>
7.1. Acyclifications . . . . .	199
7.2. Global Markov property for simple SCMs . . . . .	200
7.3. Do-calculus for simple SCMs . . . . .	201
7.4. Adjustment . . . . .	203
7.5. Bounds on causal effects . . . . .	206
<b>8. Counterfactuals</b>	<b>209</b>
8.1. Modeling counterfactuals via twinning . . . . .	209
8.2. Counterfactual Equivalence . . . . .	214
8.3. Exogenous reparameterizations . . . . .	217
8.4. Parameterizing SCMs using response functions . . . . .	220
8.5. Bounding counterfactual probabilities . . . . .	221
<b>9. Causal Discovery</b>	<b>224</b>
9.1. Detecting Causal Relations . . . . .	224
9.2. Detecting Direct Causal Relations . . . . .	226
9.3. Detecting Common Causes . . . . .	227
9.4. Abstraction through Marginalization . . . . .	230
9.5. Randomized Controlled Trials . . . . .	232
9.6. Estimating average treatment effects . . . . .	236
9.7. Faithfulness . . . . .	237
9.8. Local Causal Discovery . . . . .	238
9.9. Y-structures . . . . .	241
9.10. Minimal Separating Sets, Minimal Connecting Sets . . . . .	242
<b>10. Independence Testing</b>	<b>245</b>
10.1. Marginal Independence for Categorical Random Variables . . . . .	245
10.2. Conditional Independence for Categorical Random Variables . . . . .	251
10.3. Marginal Independence of a Random and a Non-Random Variable . . . . .	255
10.4. The general categorical case . . . . .	263

<b>11. The Fast Causal Inference Algorithm</b>	<b>266</b>
11.1. Modeling selection bias	266
11.2. Inducing walks	267
11.3. Partial Ancestral Graphs	270
11.4. Unshielded triples	278
11.5. Discriminating paths	278
11.6. Independence models and Markov equivalence	280
11.7. Skeleton search	282
11.8. FCI Algorithm	287
11.9. Completeness	293
<b>A. Appendix: Measure Theoretic Probability</b>	<b>295</b>
A.1. Why Measure Theory?	295
A.2. Core Concepts	298
A.3. Default Choices for Sigma-Algebras	300
A.4. Standard Measurable Spaces	302
A.5. Measure Integrals	303
A.6. Densities/Derivatives	306
A.7. Conditional Expectation	308
A.8. The Lebesgue Measure	310
A.9. Transformation Rules	311
A.10. Measure Extension Theorems	313
<b>References</b>	<b>318</b>

# 1. Experimental Causal Discovery

## 1.1. Types of Correlations

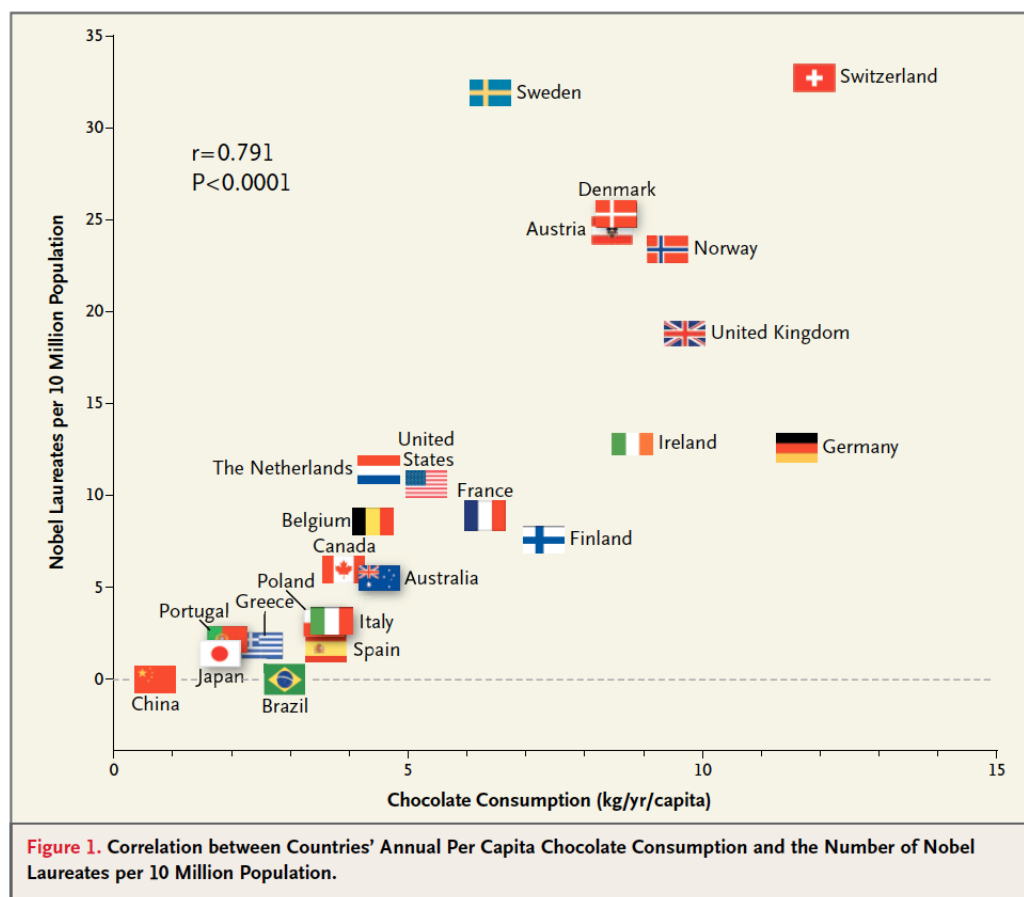


Figure 1: Correlation between chocolate consumption and Nobel prizes [Mes12].

**Explanation 1.1.1.** *What conclusions can we draw from this? Where does the correlation come from? Would the correlation hold under different conditions/circumstances? There are several explanations/stories that one could build around the measured correlation between the number of Nobel prizes  $N$  and the chocolate consumption per capita  $C$ :*

- $N$  causes  $C$ : “Nobel prize winning countries like to celebrate with chocolate consumption.”*
- $N$  is an effect of  $C$ : “Chocolate contains brain enhancing chemicals.”*
- Feedback between  $N$  and  $C$ : Both stories hold.*

- d) *Selection bias between N and C: “N and C are actually independent, but the data used was biased, e.g. only Western and Asian countries were considered. Other countries might just be in the upper left or bottom right corners.”*
- e) *Functional constraints between N and C: “International regulations make sure that Nobel prizes and chocolate imports are subtracted/added if they violate a linear relationship.”*
- f) *N and C are confounded: “The wealth of a country determines both, how much money goes to science and also how much people can spend on chocolate.”*
- g) *Other explanations, e.g. measurement error, statistical coincidence, other forms of spurious correlations, combinations of all of these, etc.?*

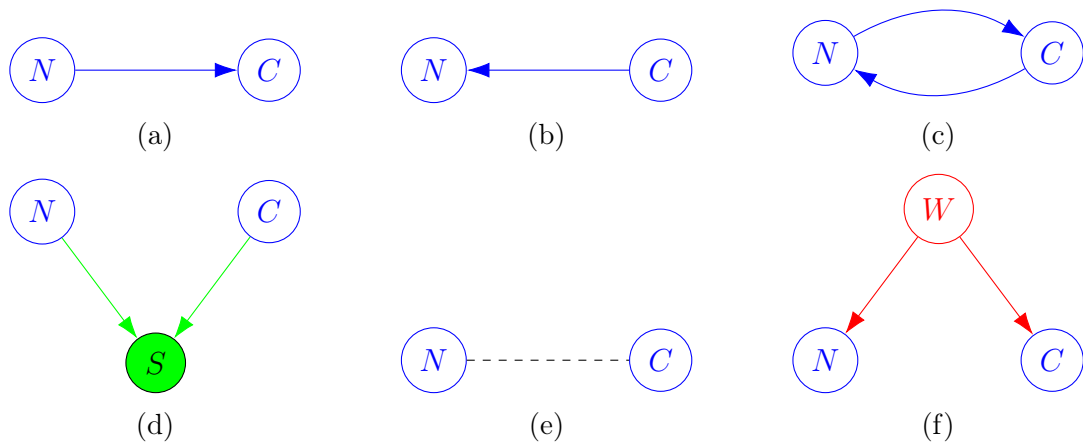


Figure 2: Graphical representations of different correlation inducing scenarios.

**Discussion 1.1.2.** *Correlation does not imply causation because there are other possible correlation inducing scenarios. Also, correlation is symmetric, causation is asymmetric.*

## 1.2. Causal Effects in the Real World

**Example 1.2.1** (Does the thermometer cause the sun to rise?). *Consider an old type of thermometer ( $T$ ) with a needle that can—for simplicity of arguments—either point to higher temperatures (up) or to lower temperatures (down). We also consider the state of the sun ( $S$ ), which can either be up ( $u$ ) or down ( $d$ ). We then observe that  $T$  correlates with  $S$ . For simplicity, we assume a one-to-one relationship:*

$T$	$S$
$u$	$u$
$d$	$d$



The conditional distribution  $P(S|T)$  then looks like this:

$$\begin{aligned} P(S = u|T = u) &= 1, \\ P(S = u|T = d) &= 0, \\ P(S = d|T = u) &= 0, \\ P(S = d|T = d) &= 1. \end{aligned}$$

If we are cold we are now tempted to try changing the needle in the thermometer in order to make the sun rise and warm us up.

What is wrong with our analysis?

**Discussion 1.2.2.** The example 1.2.1 makes clear that there is a difference between:

1. Observing the positions of the thermometer needle  $T$  and the sun  $S$ , resulting in an observational data set, leading to an estimate of  $P(S|T)$ .
2. Interacting with the thermometer needle  $T$  and getting the sun's response  $S$ , resulting in an interventional data set, leading to estimates for  $P(S|\text{do}(T))$ :

$$\begin{aligned} P(S = u|\text{do}(T = u)) &= 0.5, \\ P(S = u|\text{do}(T = d)) &= 0.5, \\ P(S = d|\text{do}(T = u)) &= 0.5, \\ P(S = d|\text{do}(T = d)) &= 0.5. \end{aligned}$$

**Definition 1.2.3** (Causal effect—real world definition). We say that a variable  $X$  has a causal effect on another variable  $Y$  if forcing  $X$  to take on a value  $x$ , the distribution of  $Y$  explicitly depends on  $x$ , that is:

$$\exists x_0, x_1 \in \mathcal{X} : \quad P(Y|\text{do}(X = x_0)) \neq P(Y|\text{do}(X = x_1)).$$

**Remark 1.2.4.** 1. Again, note that example 1.2.1 shows that the condition in definition 1.2.3 is different from:

$$\exists x_0, x_1 \in \mathcal{X} : \quad P(Y|X = x_0) \neq P(Y|X = x_1),$$

which just uses the conditional distributions instead of the interventional distributions.

2. Also note, that the 'do-operators' are not operators on the observational distribution  $P(X, Y)$  or  $P(Y|X)$ , etc., or on the corresponding observational data sets. They reflect actions/interventions in the real world leading to different distributions and corresponding data sets.
3. There are usually many possible intervention values and targets one can think of, leading to many different interventional distributions and data sets.
4. One can consider the observational distribution as a special case of an interventional distribution (where we intervene by doing nothing).

### 1.3. Randomized Controlled Trials (RCT)

**Principle 1.3.1** (Randomized Controlled Trial (RCT)). Assume we want to know if ‘treatment’ variable  $X$  has a causal effect on ‘outcome’ variable  $Y$ , i.e. we want to estimate the deviation between:  $P(Y|\text{do}(X = x_0))$  and  $P(Y|\text{do}(X = x_1))$ . For this we have test subjects  $w_1, \dots, w_N$ . A Randomized Controlled Trial then follows the following steps:

1. Split the population of test subjects into 2 groups (‘test group’  $C_1$  vs. ‘control group’  $C_0$ ) by random lot (or fair coin flips).
2. Give every test subject  $w_n \in C_1$  from ‘test group’ the treatment  $x_1$  and the ones  $w_n \in C_0$  from ‘control group’ the control treatment  $x_0$ .
3. Measure the outcome  $y_n$  for each test subject  $w_n$  and estimate the deviation:

$$D := d(P(Y|\text{do}(X = x_0)), P(Y|\text{do}(X = x_1))).$$

4. Do a statistical test if the deviation  $D$  is significantly different from 0.
5. If it is significantly different from 0 we can conclude a causal effect of  $X$  on  $Y$ , otherwise not.

**Remark 1.3.2.** The notion of a randomized controlled trial goes back several centuries. It was already described in 1648 by Flemish physician Jan Baptista van Helmont [vH48]: “Let us take from the itinerants’ hospitals, from the camps or from elsewhere 200 or 500 poor people with fevers, pleurisy etc. and divide them in two: let us cast lots so that one half of them fall to me and the other half to you. I shall cure them without blood-letting or perceptible purging, you will do so according to your knowledge (nor do I even hold you to your boast of abstaining from phlebotomy or purging) and we shall see how many funerals each of us will have: the outcome of the contest shall be the reward of 300 florins deposited by each of us. Thus shall your business be concluded. O Magistrates to whose hearts the health of your people is dear; let the trial be made for the public good, in order to know the truth, for the sake of your life and soul and for the health of all the people, sons, widows and orphans. Let there be a real debate to find the means of cure.”

**Example 1.3.3.** Example application of randomized controlled trials are:

1. drug or vaccine testing,
2. advertisement placement,
3. evaluating public policies, etc.
4. A. Banerjee, E. Duflo, M. Kremer got the Nobel Prize in Economics 2019 for using RCTs in poverty research, e.g. improving school attendance and performance in poor areas via giving different towns different incentives (e.g. text books vs. deworming medicine vs. control groups).

**Discussion 1.3.4.** 1. An RCT is an 'interventional study' (in contrast to just 'observational study') since we control the treatment and 'force' it onto the test subjects.

2. Randomized Controlled Trials are considered the gold standard for experimental causal discovery.

3. To further avoid biases one usually insists on double/triple blind RCT studies, i.e. noone directly involved in the study knows who got which treatment (e.g. neither the doctor, the experimenter, the patient, etc.).

4. Often RCTs cannot be done for ethical reasons (e.g. "smoking causes cancer" research).

5. Sometimes RCTs require too many resources to be feasible.

**Exercise 1.3.5.** Go online, find news like "drinking wine every day is good for your health" or "chewing gum causes diabetes", etc., look up the original research paper and check:

1. if they did interventional studies (like RCT) or just observational studies,

2. in case of an RCT, whether it was double/triple blind,

3. otherwise, if (and how) they ruled out other correlation inducing scenarios,

4. what bias could have possibly introduced through the data collecting process,

5. how big the data set was, what assumptions were made, what statistical methods were used, etc.,

6. what other 'stories' you could come up with in order to explain the data. Be creative, create 5 stories!

Write down your findings and talk to others about it.

## 2. Transition Probability Theory

### 2.1. Elementary Probability Theory

**Example 2.1.1** (Winning a pie with a biased die). *You are allowed to roll a biased die with 6 sides. If you roll a 5 or 6 you win a car, a 4 gives you a mug and 1, 2, 3 wins you an apple pie. In this case the sample space is  $\mathcal{W} := \{1, 2, 3, 4, 5, 6\}$  and the die introduces a probability distribution  $P$  on  $\mathcal{W}$ . Since the die is biased, we have to specify each of the probability masses to throw those numbers separately:*

$$p(1) = 0.5, \quad p(2) = p(3) = p(4) = 0.1, \quad p(5) = 0.15, \quad p(6) = 0.05.$$

*We are now interested in the probabilities of the events of winning those 3 different prizes. For this we consider the 'prize' space:  $\mathcal{Z} := \{\text{pie}, \text{mug}, \text{car}\}$ . We can then formalize the outcome via the map  $F$ :*

$$\begin{aligned} F : \quad \mathcal{W} &\rightarrow \mathcal{Z}, \\ 1, 2, 3 &\mapsto \text{pie}, \\ 4 &\mapsto \text{mug}, \\ 5, 6 &\mapsto \text{car}. \end{aligned}$$

*To compute the probability of winning each of the prizes we need to 'push' the probability distribution  $P$ , which lives on the space  $\mathcal{W}$ , to the space  $\mathcal{Z}$ . We can do this as follows:*

$$\begin{aligned} P(F = \text{pie}) &= P(F^{-1}(\{\text{pie}\})) &= P(\{1, 2, 3\}) &= p(1) + p(2) + p(3) &= 0.7, \\ P(F = \text{mug}) &= P(F^{-1}(\{\text{mug}\})) &= P(\{4\}) &= p(4) &= 0.1, \\ P(F = \text{car}) &= P(F^{-1}(\{\text{car}\})) &= P(\{5, 6\}) &= p(5) + p(6) &= 0.2, \end{aligned}$$

*where  $F^{-1}(C) := \{w \in \mathcal{W} \mid F(w) \in C\}$  is the pre-image of  $C \subseteq \mathcal{Z}$ .*

**Discussion 2.1.2.** *The simple example 2.1.1 already provides us with the main examples for the typical probability-theoretic terminology and important insights:*

- 1. We call the tuple  $(\mathcal{W}, P)$  a probability space. It is important to note that  $P$  was defined on  $\mathcal{W}$ , not  $\mathcal{Z}$ .*
- 2. We call the map  $F$  a random variable, which is really nothing else than a map from a probability space to another space.*
- 3. Events are modelled by subsets  $B \subseteq \mathcal{W}$ , not just by single elements  $w \in \mathcal{W}$ . For example consider the event that you don't win a car. This event can't be represented by a single element in  $\mathcal{W}$  or  $\mathcal{Z}$ .*
- 4. In this example we can compute the probability of an event by additivity of  $P$  and the use of the probability mass function, via  $P(B) = \sum_{w \in B} p(w)$ .*
- 5. The distribution of the prizes, i.e. the distribution of random variable  $F$ , assigns probabilities to events  $C \subseteq \mathcal{Z}$  and can be computed using the pre-image of  $F$  via  $P(F \in C) = P(F^{-1}(C))$ , where the latter is now an event  $F^{-1}(C) \subseteq \mathcal{W}$ , which we already know how to deal with.*

6. The distribution of  $F$  on  $\mathcal{Z}$  here is also called the push-forward distribution or image distribution of  $P$  via  $F$  or just the law of  $F$ . It is often abbreviated as:  $P_F$ ,  $P^F$ ,  $F_*P$  or  $P(F)$ . Again note:  $P(F)(C) := P(F \in C) = P(F^{-1}(C))$ .
7. So  $(\mathcal{Z}, P(F))$  forms a probability space on its own and as soon as we know  $P(F)$  we don't need any information about  $(\mathcal{W}, P)$  anymore if all we are interested in is the events in  $\mathcal{Z}$  and the law of  $F$ . All randomness on  $\mathcal{Z}$  is fully specified by  $P(F)$ .

**Example 2.1.3.** Now consider the standard normal distribution  $\mathcal{N}(0, 1)$  on  $\mathbb{R}$ , which is specified by the probability density function:

$$p(w) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \cdot w^2\right).$$

The probability of an event  $A \subseteq \mathbb{R}$  is then given by:

$$P(A) = \int_A p(w) dw,$$

in case  $A$  can be integrated over (i.e. if it is not a too pathological set). For instance, if  $A = [a, b] \cup [c, d]$  with  $a \leq b < c \leq d$  we get:

$$P(A) = \int_a^b p(w) dw + \int_c^d p(w) dw.$$

Note that, even though  $p(w) > 0$  for every  $w \in \mathbb{R}$ , we have:

$$P(\{x\}) = 0 \quad \text{for every } x \in \mathbb{R}.$$

Now consider the random variable  $F : \mathbb{R} \rightarrow \mathbb{R}$  with  $F(w) = \sin(w)$ . It is not immediately clear how to define the probability distribution of  $F$  when only working with probability densities. It is even more difficult to derive the probability density for  $F$  in this setting.

- Discussion 2.1.4.**
1. The examples 2.1.1 and 2.1.3 show that many probability distributions can be represented either by probability mass functions (discrete case),  $w \mapsto p(w)$ , or probability density functions (absolute continuous case),  $w \mapsto p(w)$ .
  2. Both cases have in common that one only needs a function that takes elements  $w \in \mathcal{W}$  as arguments, in contrast to subsets  $A \subseteq \mathcal{W}$ . This is usually the reason why only the discrete and absolute continuous cases are taught in elementary probability theory or machine learning classes.
  3. Note that, in the discrete case with  $K$  classes, one only needs to specify the  $K$  values  $p(1), \dots, p(K)$ , in contrast to the  $2^K$  values on subsets  $P(A)$  for  $A \in 2^{\mathcal{W}}$  (the power set of  $\mathcal{W}$  consisting of all subsets  $A \subseteq \mathcal{W}$ ), as the latter values can be derived from the former values using additivity.

4. We have problems defining probability distributions of random variables for absolute continuous distributions when we are only allowed to work with probability densities.
5. Measure theory is the framework that directly works with subsets  $A \subseteq \mathcal{W}$ , in contrast to elements  $w \in \mathcal{W}$ , and provides a unifying language that encompasses both special cases.

## 2.2. Recap - Measure Theoretic Probability

Here we just remind the reader of our notations for the core concepts of measure theoretic probability. More can be found in Appendix A.

### 2.2.1. Measurable Spaces and Maps

**Definition 2.2.1** ( $\sigma$ -algebras). Let  $\mathcal{W}$  be a set. A (non-empty) set  $\mathcal{B} \subseteq 2^{\mathcal{W}}$  of subsets  $A \subseteq \mathcal{W}$  is called a  $\sigma$ -algebra on  $\mathcal{W}$  if it satisfies the following rules:

- i) empty set:  $\emptyset \in \mathcal{B}$ ,
- ii) complement: If  $A \in \mathcal{B}$  then also:  $A^c := \mathcal{W} \setminus A \in \mathcal{B}$ ,
- iii) countable union: If  $A_n \in \mathcal{B}$  for all  $n \in \mathbb{N}$  then also:  $\bigcup_{n \in \mathbb{N}} A_n \in \mathcal{B}$ .

**Definition 2.2.2** (Measurable spaces). A tuple  $(\mathcal{W}, \mathcal{B})$  of a set  $\mathcal{W}$  and a  $\sigma$ -algebra  $\mathcal{B}$  on  $\mathcal{W}$  is called measurable space.

**Remark 2.2.3** (Abuse of notation). By abuse of notation we often just call  $\mathcal{W}$  a measurable space by implicitly assuming that it is endowed with a fixed  $\sigma$ -algebra, which we will indicate by  $\mathcal{B}_{\mathcal{W}}$  or  $\mathcal{B}(\mathcal{W})$  if needed. We will also just call a subset  $A \subseteq \mathcal{W}$  measurable when we actually mean that  $A \in \mathcal{B}_{\mathcal{W}}$ .

**Definition 2.2.4** (Measurable maps). Let  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  and  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}})$  be two measurable spaces and  $f : \mathcal{W} \rightarrow \mathcal{Z}$  be a map. We call  $f$  a  $\mathcal{B}_{\mathcal{W}}$ - $\mathcal{B}_{\mathcal{Z}}$ -measurable map (or just measurable for short) if for all  $B \in \mathcal{B}_{\mathcal{Z}}$  the pre-image  $f^{-1}(B)$  is an element of  $\mathcal{B}_{\mathcal{W}}$ . In formulas:

$$\forall B \in \mathcal{B}_{\mathcal{Z}} : f^{-1}(B) \in \mathcal{B}_{\mathcal{W}}.$$

Remember the definition of pre-image:  $f^{-1}(B) := \{w \in \mathcal{W} \mid f(w) \in B\}$ .

For most of the lecture we will restrict to well-behaved measurable spaces, namely standard measurable spaces. The key point is that they all behave like the space  $[0, 1]$ , or  $\mathbb{R}$ , with its Borel- $\sigma$ -algebra. So (almost) all results for  $[0, 1]$  immediately translate to standard measurable spaces.

**Definition 2.2.5** (Standard measurable space, see [Fre15] 424A-G). A measurable space  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  is called standard measurable space (aka standard Borel space) if it is measurably isomorphic to either:

1. a finite measurable space  $\{1, \dots, M\}$  for some  $M \in \mathbb{N}$  endowed with the power set  $\sigma$ -algebra  $2^{\{1, \dots, M\}}$ , or:
2. the countably infinite space  $\mathbb{N}$  endowed with the power set  $\sigma$ -algebra  $2^{\mathbb{N}}$ , or:
3. the unit interval  $[0, 1]$  endowed with its Borel  $\sigma$ -algebra<sup>1</sup>:

$$\mathcal{B}_{[0,1]} = \sigma(\{[a, b] \mid a, b \in [0, 1] \cap \mathbb{Q}, a \leq b\}).$$

“Measurably isomorphic” means that there is a measurable map from one space to the other that has a measurable inverse.

The following theorem shows that (almost) all spaces we encounter in practice are actually standard measurable spaces, justifying our focus on standard measurable spaces for the most of this lecture.

**Theorem 2.2.6** (Kuratowski et al., see [Fre15] 424A-G). *Every Borel subset of any complete metric space that has a countable dense subset is a standard measurable space in its Borel  $\sigma$ -algebra.*

**Example 2.2.7.**  $\mathbb{R}, \mathbb{R}^D, \mathbb{Q}, \mathbb{Z}, \mathbb{N}, \{1, \dots, M\}, [0, 1]$ , topological manifolds, countable CW-complexes, etc., are all standard measurable spaces.

## 2.2.2. Finite and Probability Measures

**Definition 2.2.8** (Measures). *Let  $(\mathcal{W}, \mathcal{B})$  be a measurable space. A measure  $\mu$  on  $(\mathcal{W}, \mathcal{B})$ —by definition—is a map:*

$$\mu : \mathcal{B} \rightarrow \mathbb{R} \cup \{\infty\}, \quad D \mapsto \mu(D),$$

such that:

- i) non-negative:  $\forall A \in \mathcal{B}: \mu(A) \in [0, \infty]$ ,
- ii) empty set:  $\mu(\emptyset) = 0$ ,
- iii) countably additive (aka  $\sigma$ -additive): for all sequences  $A_n \in \mathcal{B}$ ,  $n \in \mathbb{N}$ , with  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ , we have:

$$\mu\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} \mu(A_n).$$

**Definition 2.2.9** (Probability and finite measures). *A measure  $\mu$  on  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  is called:*

1. probability measure if  $\mu(\mathcal{W}) = 1$ .

---

<sup>1</sup>See Definition A.3.2 for the  $\sigma$ -algebra generated by a set of subsets.

2. finite measure if  $\mu(\mathcal{W}) < \infty$ .

3.  $\sigma$ -finite measure if there are  $D_n \in \mathcal{B}$ ,  $n \in \mathbb{N}$ , with  $\mu(D_n) < \infty$  and  $\mathcal{W} = \bigcup_{n \in \mathbb{N}} D_n$ .

Clearly, every probability measure is finite, and, every finite measure is  $\sigma$ -finite.

**Definition 2.2.10** (The spaces of finite and probability measures). *The set of all probability measures on  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  is denoted by  $\mathcal{P}(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$ , and the set of all finite measures by  $\mathcal{M}(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$ , or  $\mathcal{P}(\mathcal{W})$  and  $\mathcal{M}(\mathcal{W})$ , resp., for short. For  $B \in \mathcal{B}_{\mathcal{W}}$  we consider the evaluation map:*

$$\text{ev}_B : \mathcal{M}(\mathcal{W}) \rightarrow \mathbb{R}_{\geq 0}, \quad \mu \mapsto \text{ev}_B(\mu) := \mu(B).$$

We then endow  $\mathcal{M}(\mathcal{W})$ , and  $\mathcal{P}(\mathcal{W})$ , resp., with the smallest  $\sigma$ -algebra  $\mathcal{B}$  such that all evaluation maps  $\text{ev}_B$  are  $\mathcal{B}$ - $\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable, where  $\mathcal{B}_{\mathbb{R}_{\geq 0}}$  is the Borel- $\sigma$ -algebra of  $\mathbb{R}_{\geq 0}$ , i.e.:

$$\mathcal{B}_{\mathcal{M}(\mathcal{W})} := \sigma \left( \left\{ \text{ev}_B^{-1}((r, \infty)) \mid B \in \mathcal{B}, r \in \mathbb{R}_{\geq 0} \right\} \right).$$

**Remark 2.2.11.** *The above definition implies that for measurable spaces  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ ,  $(\mathcal{Y}, \mathcal{B}_{\mathcal{Y}})$ , a map:*

$$K : \mathcal{X} \rightarrow \mathcal{M}(\mathcal{Y}),$$

is  $\mathcal{B}_{\mathcal{X}}$ - $\mathcal{B}_{\mathcal{M}(\mathcal{Y})}$ -measurable if and only if for all  $B \in \mathcal{B}_{\mathcal{Y}}$  the composition:

$$\text{ev}_B \circ K : \mathcal{X} \rightarrow \mathcal{M}(\mathcal{Y}) \rightarrow \mathbb{R}_{\geq 0},$$

is  $\mathcal{B}_{\mathcal{X}}$ - $\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable. Similarly, for  $\mathcal{P}(\mathcal{Y})$ .

**Theorem 2.2.12** (See [Par05] Thm. 6.2 + 6.5 or [Fre15] 437R). *If  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  is a standard measurable space then also  $\mathcal{P}(\mathcal{W})$  is a standard measurable space (in its usual  $\sigma$ -algebra).*

### 2.2.3. The Measure Integral

For a measure  $\mu$  on a measurable space  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$  the measure integral of measurable functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  is treated in Appendix A.5. Here we just want to remind the reader of our several different notations, which we will use interchangeably during the course:

**Notation 2.2.13** (Measure integral). *We abbreviate the measure integral of a measurable function  $f : \mathcal{X} \rightarrow \mathbb{R}$  w.r.t. measure  $\mu$  on  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$  as:*

$$\int f d\mu = \int f(x) d\mu(x) = \int f(x) \mu(dx).$$

If  $P$  is a probability measure on  $\mathcal{X} = \mathbb{R}^D$  that is either discrete or absolute continuous we have:

$$\int f(x) P(dx) = \begin{cases} \sum_{x \in \mathcal{X}} f(x) \cdot p(x), & \text{if } P \text{ is discrete,} \\ \int_{\mathcal{X}} f(x) \cdot p(x) dx & \text{if } P \text{ is absolute continuous,} \end{cases}$$

where  $p$  either denotes the probability mass function or the probability density, resp.



## 2.2.4. The Lebesgue Measure

By far the most important measure is the Lebesgue measure, which assigns the typical  $D$ -dimensional volume to cubes, i.e. the product of their side lengths.

**Definition 2.2.14** (The Lebesgue (outer) measure). *The Lebesgue (outer) measure  $\lambda^D$  on  $\mathbb{R}^D$  is given for subsets  $A \subseteq \mathbb{R}^D$  via:*

$$\lambda^D(A) := \inf \left\{ \sum_{n \in \mathbb{N}} \text{vol}^D([a^{(n)}, b^{(n)}]) \mid A \subseteq \bigcup_{n \in \mathbb{N}} [a^{(n)}, b^{(n)}] \right\},$$

where the infimum is running over sequences of  $D$ -dimensional cubes:

$$[a^{(n)}, b^{(n)}] = [a_1^{(n)}, b_1^{(n)}] \times \cdots \times [a_D^{(n)}, b_D^{(n)}],$$

with  $a^{(n)} = (a_1^{(n)}, \dots, a_D^{(n)})$ ,  $b^{(n)} = (b_1^{(n)}, \dots, b_D^{(n)}) \in \mathbb{R}^D$ ,  $a_d^{(n)} \leq b_d^{(n)}$  for  $d = 1, \dots, D$ ,  $n \in \mathbb{N}$ , that jointly cover  $A$ , where the  $D$ -dimensional volume is given by:

$$\text{vol}^D([a^{(n)}, b^{(n)}]) := (b_1^{(n)} - a_1^{(n)}) \cdots (b_D^{(n)} - a_D^{(n)}), \quad \text{vol}^D(\emptyset) := 0.$$

**Theorem 2.2.15** (The Lebesgue measure). *The Lebesgue measure  $\lambda^D$ , when restricted to the Borel- $\sigma$ -algebra of  $\mathbb{R}^D$ , is the unique measure on  $\mathbb{R}^D$  that satisfies:*

$$\lambda^D([a, b]) = \text{vol}^D([a, b]),$$

for all  $D$ -dimensional cubes  $[a, b]$ . If the dimension is clear from the context we might just write  $\lambda$  for  $\lambda^D$ .

## 2.3. Transition Measures and Markov Kernels

### 2.3.1. Core Definitions

**Motivation 2.3.1.** *If we consider a deterministic measurable map  $f : \mathcal{T} \rightarrow \mathcal{W}$  then  $f$  assigns to each point  $t \in \mathcal{T}$  exactly one point  $w = f(t) \in \mathcal{W}$ . Sometimes we rather want to model a probabilistic map, i.e. an assignment that can be random or comes with some uncertainties but still changes depending on the input  $t$ . The notion of Markov kernels formalizes this. A Markov kernel  $K$  from  $\mathcal{T}$  to  $\mathcal{W}$  can be considered a measurable map from  $\mathcal{T}$  to the space of probability measures  $\mathcal{P}(\mathcal{W})$  of  $\mathcal{W}$ :*

$$\mathcal{T} \rightarrow \mathcal{P}(\mathcal{W}).$$

*It assigns to each  $t \in \mathcal{T}$  a probability distribution over  $\mathcal{W}$ , which then assigns to each measurable subset  $D \subseteq \mathcal{W}$  a probability value in  $[0, 1]$ .*

**Example 2.3.2** (Markov kernels). *1. any statistical model (i.e. family of model distributions)  $\{p_\theta \mid \theta \in \mathcal{F}\}$ , can be considered a Markov kernel, which we write  $P(X|\Theta)$ .*

*2. any conditional distribution  $P(Y|X)$  can be considered a Markov kernel.*

3. a neural network with softmax output for classification with input  $x \in \mathcal{X}$ , output  $y \in \mathcal{Y}$  and weights  $w \in \mathcal{W}$  can be seen as a Markov kernel  $P(Y|X, W)$ .

We first start slightly more generally by defining finite transition measures.

**Definition 2.3.3** (Finite transition measures and Markov kernels). *Let  $\mathcal{T}$ ,  $\mathcal{W}$  be measurable spaces.*

1. A (finite<sup>2</sup>) transition measure from  $\mathcal{T}$  to  $\mathcal{W}$  is—per definition—a measurable map:

$$K : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

from  $\mathcal{T}$  to the space of finite measures of  $\mathcal{W}$ .

2. A transition probability or Markov kernel is—per definition—a measurable map:

$$K : \mathcal{T} \rightarrow \mathcal{P}(\mathcal{W}),$$

from  $\mathcal{T}$  to the space of probability measures of  $\mathcal{W}$ .

**Notation 2.3.4** (Transition measures and Markov kernels). *1. We often use suggestive notations as follows for finite transition measures and Markov kernels:*

$$K(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}), \quad t \mapsto K(W|T = t),$$

where for every fixed  $t \in \mathcal{T}$  the following map:

$$K(W|T = t) : \mathcal{B}_{\mathcal{W}} \rightarrow \mathbb{R}_{\geq 0}, \quad D \mapsto K(W \in D|T = t),$$

is a finite measure, or probability measure, respectively.

2. For fixed  $D \in \mathcal{B}_{\mathcal{W}}$  we then use the following notation for the following measurable map:

$$K(W \in D|T) : \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}, \quad t \mapsto K(W \in D|T = t).$$

3. Since  $K(W|T)$  takes the argument  $t \in \mathcal{T}$  first, but then also  $D \in \mathcal{B}_{\mathcal{W}}$  as a second argument we can also indentify  $K(W|T)$  with the following two-argument map, which we denote with the same symbols:

$$K(W|T) : \mathcal{B}_{\mathcal{W}} \times \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}, \quad (D, t) \mapsto K(W \in D|T = t).$$

4. For Markov kernels  $K(W|T)$  we will most of the time use the dashed arrow to  $\mathcal{W}$  (instead of a usual arrow to  $\mathcal{P}(\mathcal{W})$ ) to indicate the Markov kernel as follows:

$$K(W|T) : \mathcal{T} \dashrightarrow \mathcal{W}, \quad (D, t) \mapsto K(W \in D|T = t).$$

---

<sup>2</sup>In this course we will only discuss *finite* transition measures and just drop the word “finite” for simplicity in the following.

5. Note that above  $W$  and  $T$  are considered suggestive symbols only, but one could give  $W$  the meaning to mean the (identity or) projection map  $\text{pr}_{\mathcal{W}}$  onto  $\mathcal{W}$ . From the point on we also have a map  $T$  mapping to  $\mathcal{T}$  the notation becomes ambiguous:  $K(W|T)$  could also mean  $K(W|T)$  where we plugged in  $T$  for  $t$  in “ $T = t$ ”, similar to conditional expectations  $\mathbb{E}[W|T]$ , but the meaning should become clear from the context.

The implicit correspondence in the above discussion can more formally be summarized as:

**Lemma 2.3.5.** *There is a one-to-one correspondence between the following constructions:*

1. a finite transition measure, i.e. a measurable map:

$$K(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}), \quad t \mapsto K(W|T = t).$$

2. a two-argument function:

$$\tilde{K}(W|T) : \mathcal{B}_{\mathcal{W}} \times \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}, \quad (D, t) \mapsto \tilde{K}(W \in D|T = t),$$

such that:

- i) For each  $t \in \mathcal{T}$  the map:

$$\mathcal{B}_{\mathcal{W}} \rightarrow \mathbb{R}_{\geq 0}, \quad D \mapsto \tilde{K}(W \in D|T = t)$$

is a finite measure (i.e. countably additive with  $\tilde{K}(W \in \mathcal{W}|T = t) < \infty$  for all  $t \in \mathcal{T}$ ).<sup>3</sup>

- ii) For each  $D \in \mathcal{B}_{\mathcal{W}}$  the map:

$$\mathcal{T} \rightarrow \mathbb{R}_{\geq 0}, \quad t \mapsto \tilde{K}(W \in D|T = t)$$

is  $\mathcal{B}_{\mathcal{T}}\text{-}\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable.

For Markov kernels the same statement holds after replacing  $\mathcal{M}(\mathcal{W})$  with  $\mathcal{P}(\mathcal{W})$  and “finite measure” with “probability measure”.

*Proof.* The correspondence is via putting  $K(W \in D|T = t) = \tilde{K}(W \in D|T = t)$  and vice versa. The corresponding properties hold by definition of the  $\sigma$ -algebra on  $\mathcal{M}(\mathcal{W})$ , also see Remark 2.2.11. Working out the details is left as an exercise.  $\square$

---

<sup>3</sup>Note that for a finite transition measure the finite value  $\tilde{K}(W \in \mathcal{W}|T = t)$  can vary with  $t \in \mathcal{T}$ . This is in contrast to Markov kernels where we always have  $\tilde{K}(W \in \mathcal{W}|T = t) = 1$  for all  $t \in \mathcal{T}$ .

### 2.3.2. Special Cases of Markov Kernels

**Example 2.3.6** (Markov kernels on discrete spaces). *Consider a Markov kernel:*

$$K(W|T) : \mathcal{T} \dashrightarrow \mathcal{W}, \quad (D, t) \mapsto K(W \in D|T = t),$$

where both  $\mathcal{W} = \{w_1, \dots, w_M\}$  and  $\mathcal{T} = \{t_1, \dots, t_K\}$  are finite discrete spaces. Then we can define the mass function  $k$  via:

$$k(w_i|t_j) := K(W \in \{w_i\}|T = t_j),$$

and the matrix  $\tilde{K} := (k(w_i|t_j))_{i,j}$ . Then the matrix  $\tilde{K}$  is a stochastic matrix, i.e. it has non-negative entries and each of its columns sums to 1.  $\tilde{K}$  then fully determines the Markov kernel  $K$ . So in the (finite) discrete case a Markov kernel is basically nothing else than a stochastic matrix filled with the transition probabilities.

**Example 2.3.7** (Linear Gaussian Markov kernels). *Let  $\mathcal{W} = \mathbb{R}^M$ ,  $\mathcal{T} = \mathbb{R}^L$ ,  $\gamma \in \mathbb{R}^M$ ,  $\Gamma \in \mathbb{R}^{M \times L}$  and  $\Sigma \in \mathbb{R}^{M \times M}$  a fixed symmetric, positive-definite covariance matrix. Then:*

$$K(W \in D|T = t) := \int_D \mathcal{N}(w|\Gamma \cdot t + \gamma, \Sigma) dw,$$

defines a Markov kernel from  $\mathcal{T}$  to  $\mathcal{W}$ . Markov kernels of this form are called linear Gaussian Markov kernels. If  $\Sigma$  is only positive-semi-definite we call  $K(W|T)$  a degenerate or generalized linear Gaussian Markov kernel.

**Example 2.3.8** (Exponential families as finite transition measures). *Let  $\mathcal{W}$  be a measurable space and  $\mu$  a (non-zero) measure on  $\mathcal{W}$  and  $S : \mathcal{W} \rightarrow \mathbb{R}^D$  a measurable map. Define for  $t \in \mathbb{R}^D$ :*

$$Z(t) := \int_{\mathcal{W}} \exp(t^\top S(w)) \mu(dw) \in (0, \infty].$$

We then put:

$$\mathcal{T} := \{t \in \mathbb{R}^D \mid Z(t) < \infty\}.$$

We can then define the finite transition measure  $K(W|T)$  from  $\mathcal{T}$  to  $\mathcal{W}$  for  $D \in \mathcal{B}_{\mathcal{W}}$  and  $t \in \mathcal{T}$  via:

$$K(W \in D|T = t) := \int_D \exp(t^\top S(w)) \mu(dw).$$

From this we get the Markov kernel  $Q(W|T) : \mathcal{T} \dashrightarrow \mathcal{W}$  via normalization:

$$Q(W \in D|T = t) := \int_D \exp(t^\top S(w) - L(t)) \mu(dw),$$

with log-normalizer:  $L(t) := \log Z(t) = \log K(W \in \mathcal{W}|T = t)$ .

**Remark 2.3.9** (Markov kernels generalize probability distributions). *Let  $\mathcal{W}$  be a measurable space.*

1. Every probability distribution  $P \in \mathcal{P}(\mathcal{W})$  can be considered as a constant Markov kernel from  $\mathcal{T}$  to  $\mathcal{W}$  via:

$$K : \mathcal{T} \dashrightarrow \mathcal{W}, \quad (D, t) \mapsto K(D|t) := P(D).$$

2. Every Markov kernel from the one-point space:  $\mathcal{T} = * := \{*\}$  to  $\mathcal{W}$ :

$$K : * \dashrightarrow \mathcal{W}, \quad (D, *) \mapsto K(D|*),$$

defines a unique probability distribution  $P \in \mathcal{P}(\mathcal{W})$  given via:

$$P(D) := K(D|*).$$

So we can identify probability distributions on  $\mathcal{W}$  with Markov kernels  $* \dashrightarrow \mathcal{W}$ .

**Remark 2.3.10** (Markov kernels generalize deterministic maps). Consider a measurable mapping  $f : \mathcal{T} \rightarrow \mathcal{W}$ . Then we can turn  $f$  into a Markov kernel  $\delta_f$  via:

$$\delta_f : \mathcal{T} \dashrightarrow \mathcal{W}, \quad (D, t) \mapsto \delta_f(D|t) := \mathbb{1}_D(f(t)),$$

which puts 100% probability mass onto the function value  $f(t)$  for given  $t \in \mathcal{T}$ .

### 2.3.3. The Doob-Radon-Nikodym Derivative

**Definition 2.3.11** (Absolute continuity). Let  $\mathcal{T}, \mathcal{W}$  be measurable spaces and

$$Q(W|T), K(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

two finite transition measures. We say that  $Q(W|T)$  is absolute continuous w.r.t.  $K(W|T)$  if for all  $t \in \mathcal{T}$  and  $D \in \mathcal{B}_{\mathcal{W}}$  we have the implication:

$$K(W \in D|T = t) = 0 \quad \implies \quad Q(W \in D|T = t) = 0.$$

In symbols we abbreviate this as:

$$Q(W|T) \ll K(W|T).$$

**Remark 2.3.12.** For absolute continuous finite transition measures  $Q(W|T) \ll K(W|T)$  there exists by the Theorem of Radon-Nikodym, see Theorem A.6.4 or [Kle20] Cor. 7.34, for each  $t \in \mathcal{T}$  separately a Radon-Nikodym derivative, i.e. a  $\mathcal{B}_{\mathcal{W}}\text{-}\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable map:

$$g_t : \mathcal{W} \rightarrow \mathbb{R}_{\geq 0},$$

such that for all  $D \in \mathcal{B}_{\mathcal{W}}$ :

$$Q(W \in D|T = t) = \int \mathbb{1}_D(w) \cdot g_t(w) K(W \in dw|T = t).$$

Unfortunately, the map:

$$g : \mathcal{W} \times \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}, \quad g(w|t) := g_t(w),$$

is not guaranteed to be jointly measurable, i.e.  $(\mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}})$ - $\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable. In case it was, we would call it a Doob-Radon-Nikodym derivative of  $Q(W|T)$  w.r.t.  $K(W|T)$ . Doob invented an alternative, but a bit more restrictive approach than the usual one to construct Radon-Nikodym derivatives for measures based on martingales. This approach will be seen to also work for the construction of Doob-Radon-Nikodym derivatives for finite transition measures.

**Definition 2.3.13** (Doob-Radon-Nikodym derivative). *Let  $\mathcal{T}, \mathcal{W}$  be measurable spaces and*

$$Q(W|T), K(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

*two finite transition measures. A map*

$$g : \mathcal{W} \times \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}, \quad (w, t) \mapsto g(w|t),$$

*is called Doob-Radon-Nikodym derivative if  $g$  is  $(\mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}})$ - $\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable and for all  $t \in \mathcal{T}$  and all  $D \in \mathcal{B}_{\mathcal{W}}$  we have:*

$$Q(W \in D|T = t) = \int \mathbb{1}_D(w) \cdot g(w|t) K(W \in dw|T = t).$$

*In other words,  $g$  provides a Radon-Nikodym derivative simultaneously for all  $t \in \mathcal{T}$ :*

$$g(w|t) = \frac{Q(W \in dw|T = t)}{K(W \in dw|T = t)}(w),$$

*that is even jointly measurable in  $(w, t)$ .*

**Lemma 2.3.14.** *If  $Q(W|T)$  has a Doob-Radon-Nikodym derivative w.r.t.  $K(W|T)$  then  $Q(W|T)$  is absolute continuous w.r.t.  $K(W|T)$ .*

*Proof.* Should be clear, left as an exercise. □

To investigate the uniqueness of the Doob-Radon-Nikodym derivative we need the following notion of  $K(W|T)$ -null sets.

**Definition 2.3.15** (Null sets). *Let  $K(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W})$  be a finite transition measure. A subset  $N \subseteq \mathcal{W} \times \mathcal{T}$  is called  $K(W|T)$ -null if  $N_t := \{w \in \mathcal{W} \mid (w, t) \in N\}$  is a  $K(W|T = t)$ -null set for every  $t \in \mathcal{T}$ , i.e. if for every  $t \in \mathcal{T}$  there exists a measurable set  $M_t \in \mathcal{B}_{\mathcal{W}}$  such that  $K(W \in M_t|T = t) = 0$  and  $N_t \subseteq M_t$ .*

**Lemma 2.3.16** (Essential uniqueness of the Doob-Radon-Nikodym derivative). *Let  $\mathcal{T}, \mathcal{W}$  be measurable spaces and:*

$$Q(W|T), K(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

be two finite transition measures with  $Q(W|T) \ll K(W|T)$  and let  $g_1, g_2$  be two Doob-Radon-Nikodym derivatives. Then the set:

$$N := \{(w, t) \in \mathcal{W} \times \mathcal{T} \mid g_1(w|t) \neq g_2(w|t)\}$$

is a  $K(W|T)$ -null set and an element of the product  $\sigma$ -algebra  $\mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}}$ . In this sense, the Doob-Radon-Nikodym derivative is essentially unique.

**Theorem 2.3.17** (Doob-Radon-Nikodym, see [DM83] Thm. 58, [Kle20] Ex. 11.17). *Let  $\mathcal{T}, \mathcal{W}$  be measurable spaces and:*

$$K(W|T), Q(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

be two finite transition measures. Assume that  $\mathcal{W}$  is a standard measurable space.<sup>4</sup> Then the following two statements are equivalent:

1.  $Q(W|T)$  is absolute continuous w.r.t.  $K(W|T)$ .
2.  $Q(W|T)$  has a Doob-Radon-Nikodym derivative w.r.t.  $K(W|T)$ .

In that case the Doob-Radon-Nikodym derivative is essentially unique.

**Remark 2.3.18.** 1. As mentioned in the footnote<sup>4</sup> Theorem 2.3.17 still holds if one only requires  $\mathcal{B}_{\mathcal{W}}$  to be countably generated. Further extensions could be made to  $\sigma$ -algebras  $\mathcal{B}_{\mathcal{W}}$  that are countably generated up to some form of null-sets.

2. With more technical conditions one could extend Theorem 2.3.17 to work for  $\sigma$ -finite transition measures. A simple, but important, special case is treated in the following Corollary 2.3.19.

**Corollary 2.3.19** (Doob-Radon-Nikodym derivatives w.r.t.  $\sigma$ -finite measures). *Let  $\mathcal{T}, \mathcal{W}$  be measurable spaces, where  $\mathcal{W}$  is a standard<sup>4</sup> measurable space, let*

$$P(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

be a finite transition measure and  $\mu$  be a  $\sigma$ -finite measure on  $\mathcal{W}$ . Then the following two statements are equivalent:

1.  $P(W|T)$  is absolute continuous w.r.t.  $\mu$ , i.e. for all  $t \in \mathcal{T}$  and  $D \in \mathcal{B}_{\mathcal{W}}$ :

$$\mu(D) = 0 \quad \implies \quad P(W \in D|T = t) = 0.$$

2.  $P(W|T)$  has a Doob-Radon-Nikodym derivative w.r.t.  $\mu$ , i.e. a jointly measurable map:

$$p : \mathcal{W} \times \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}, \quad (w, t) \mapsto p(w|t),$$

such that for all  $t \in \mathcal{T}$  and  $D \in \mathcal{B}_{\mathcal{W}}$ :

$$P(W \in D|T = t) = \int_D p(w|t) \mu(dw).$$

---

<sup>4</sup>The proof shows that we actually only require that  $\mathcal{B}_{\mathcal{W}}$  is countably generated.

In that case the Doob-Radon-Nikodym derivative is essentially unique, i.e. for two such  $p$ , say  $p_1$  and  $p_2$ , the set:

$$N := \{(w, t) \in \mathcal{W} \times \mathcal{T} \mid p_1(w|t) \neq p_2(w|t)\} \in \mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}},$$

satisfies  $\mu(N_t) = 0$  for all  $t \in \mathcal{T}$ .

*Proof.* Let  $P(W|T)$  be absolute continuous w.r.t.  $\mu$ . Since  $\mu$  is  $\sigma$ -finite there exists a probability measure  $Q(W)$  with:

$$Q(W) \ll \mu \ll Q(W).$$

Indeed, if  $\mu$  is finite, we can just put  $Q(W \in D) := \frac{\mu(D)}{\mu(\mathcal{W})}$ . If  $\mu$  is  $\sigma$ -finite, but not finite, then we have a decomposition  $\mathcal{W} = \dot{\bigcup}_{n \in \mathbb{N}} \mathcal{W}_n$  with  $0 < \mu(\mathcal{W}_n) < \infty$ . We can then put:

$$Q(W \in D) := \sum_{n \in \mathbb{N}} 2^{-n} \frac{\mu(D \cap \mathcal{W}_n)}{\mu(\mathcal{W}_n)}.$$

By the standard Radon-Nikodym theorem there exists a Radon-Nikodym derivative  $q$  of  $Q(W)$  w.r.t.  $\mu$ . Note that  $Q(W)$  defines the constant Markov kernel  $Q(W|T)$  via  $Q(W|T = t) := Q(W)$ . We thus have the absolute continuity:

$$P(W|T) \ll \mu \ll Q(W|T).$$

By Theorem 2.3.17 we thus get a Doob-Radon-Nikodym derivative  $k$  of  $P(W|T)$  w.r.t.  $Q(W|T)$ . Then  $p$  given by:

$$p(w|t) := k(w|t) \cdot q(w),$$

is a Doob-Radon-Nikodym derivative of  $P(W|T)$  w.r.t.  $\mu$ . Indeed, we get for all  $t \in \mathcal{T}$  and  $D \in \mathcal{B}_{\mathcal{W}}$ :

$$P(W \in D|T = t) = \int_D k(w|t) Q(W \in dw|T = t) = \int_D k(w|t) \cdot q(w) \mu(dw).$$

This shows one direction.

The essential uniqueness follows similar to Lemma 2.3.16 and the other direction similar to Lemma 2.3.14.  $\square$

**Corollary 2.3.20** (Absolute continuity and strictly positive densities). *Let  $\mathcal{T}$ ,  $\mathcal{W}$  be measurable spaces, where  $\mathcal{W}$  is a standard<sup>4</sup>, and:*

$$P(W|T), K(W|T), Q(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

be finite transition measures and  $\mu$  be a  $\sigma$ -finite measure on  $\mathcal{W}$ .

1.  $Q(W|T)$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $K(W|T)$  if and only if:

$$Q(W|T) \ll K(W|T) \ll Q(W|T).$$



2.  $P(W|T)$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $\mu$  if and only if:

$$\mu \ll P(W|T) \ll \mu.$$

*Proof.* The second case follows from the first using the arguments from Corollary 2.3.19. So, first assume that  $Q(W|T)$  has a strictly positive density  $q > 0$  w.r.t.  $K(W|T)$ . Then by Lemma 2.3.14 we already have:  $Q(W|T) \ll K(W|T)$ . Since  $q$  is strictly positive we can put for  $w \in \mathcal{W}$  and  $t \in \mathcal{T}$ :

$$k(w|t) := \frac{1}{q(w|t)} > 0.$$

Then  $k$  is a (strictly positive) density of  $K(W|T)$  w.r.t.  $Q(W|T)$  and we again can use Lemma 2.3.14 to also get:  $K(W|T) \ll Q(W|T)$ . This shows one direction.

Now assume that we have:

$$Q(W|T) \ll K(W|T) \ll Q(W|T).$$

Then by the Doob-Radon-Nikodym Theorem 2.3.17 we have Doob-Radon-Nikodym derivatives  $q$  and  $k$  of  $Q(W|T)$  w.r.t.  $K(W|T)$  and of  $K(W|T)$  w.r.t.  $Q(W|T)$ , resp. For all  $D \in \mathcal{B}_{\mathcal{W}}$  and  $t \in \mathcal{T}$  we thus get:

$$\begin{aligned} \int_D 1 Q(W \in dw|T = t) &= Q(W \in D|T = t) \\ &= \int_D q(w|t) K(W \in dw|T = t) \\ &= \int_D q(w|t) \cdot k(w|t) Q(W \in dw|T = t). \end{aligned}$$

Since this holds for all  $D \in \mathcal{B}_{\mathcal{W}}$  and  $t \in \mathcal{T}$  we get that the set:

$$N := \{(w, t) \in \mathcal{W} \times \mathcal{T} \mid 1 \neq q(w|t) \cdot k(w|t)\} \in \mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}},$$

is a  $Q(W|T)$ -null set, and because  $K(W|T) \ll Q(W|T)$ , also a  $K(W|T)$ -null set. We then put:

$$\tilde{q}(w|t) := q(w|t) \cdot \mathbb{1}_{N^c}(w, t) + \mathbb{1}_N(w, t), \quad (1)$$

$$\tilde{k}(w|t) := k(w|t) \cdot \mathbb{1}_{N^c}(w, t) + \mathbb{1}_N(w, t). \quad (2)$$

These are then still corresponding Doob-Radon-Nikodym derivatives and satisfy for all  $w \in \mathcal{W}$  and  $t \in \mathcal{T}$ :

$$\tilde{q}(w|t) \cdot \tilde{k}(w|t) = 1,$$

which directly implies:  $\tilde{q}(w|t), \tilde{k}(w|t) > 0$  for all  $w \in \mathcal{W}$  and  $t \in \mathcal{T}$ . This shows the claim.  $\square$

## Proofs - Theorem of Doob-Radon-Nikodym

**Lemma 2.3.21** (Essential uniqueness of the Doob-Radon-Nikodym derivative). *Let  $\mathcal{T}$ ,  $\mathcal{W}$  be measurable spaces and:*

$$Q(W|T), K(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

*be two finite transition measures with  $Q(W|T) \ll K(W|T)$  and let  $g_1, g_2$  be two Doob-Radon-Nikodym derivatives. Then the set:*

$$N := \{(w, t) \in \mathcal{W} \times \mathcal{T} \mid g_1(w|t) \neq g_2(w|t)\}$$

*is a  $K(W|T)$ -null set and an element of the product  $\sigma$ -algebra  $\mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}}$ .*

*Proof.* Consider the set:

$$N^> := \{(w, t) \in \mathcal{W} \times \mathcal{T} \mid g_1(w|t) > g_2(w|t)\} = (g_1 \times g_2)^{-1}(\Delta^>),$$

where  $\Delta^>$  is the measurable set:

$$\Delta^> := \{(r_1, r_2) \in \mathbb{R} \times \mathbb{R} \mid r_1 > r_2\} \in \mathcal{B}_{\mathbb{R}^2}.$$

Since both  $g_1$  and  $g_2$  are jointly measurable that shows that  $N^> \in \mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}}$ . It follows that  $N_t^> \in \mathcal{B}_{\mathcal{W}}$ . Furthermore, we get:

$$\begin{aligned} 0 &= Q(W \in N_t^> | T = t) - Q(W \in N_t^> | T = t) \\ &= \int \mathbb{1}_{N_t^>}(w) \cdot g_1(w|t) K(W \in dw | T = t) - \int \mathbb{1}_{N_t^>}(w) \cdot g_2(w|t) K(W \in dw | T = t) \\ &= \int \underbrace{\mathbb{1}_{N_t^>}(w) \cdot (g_1(w|t) - g_2(w|t))}_{>0 \text{ for } w \in N_t^>} K(W \in dw | T = t). \end{aligned}$$

This shows that  $K(W \in N_t^> | T = t) = 0$ . By flipping  $g_1$  and  $g_2$  we also get:  $K(W \in N_t^< | T = t) = 0$  and thus  $K(W \in N_t | T = t) = 0$ , where we notice that  $N = N^> \cup N^<$ . This shows the claim.  $\square$

**Theorem 2.3.22** (Existence of the Doob-Radon-Nikodym derivative, see [DM83] Thm. 58, [Kle20] Ex. 11.17). *Let  $\mathcal{T}$ ,  $\mathcal{W}$  be measurable spaces and:*

$$K(W|T), Q(W|T) : \mathcal{T} \rightarrow \mathcal{M}(\mathcal{W}),$$

*be two finite transition measures. Assume that  $\mathcal{W}$  is a standard measurable space.<sup>4</sup>  $Q(W|T) \ll K(W|T)$  implies that  $Q(W|T)$  has a Doob-Radon-Nikodym derivative w.r.t.  $K(W|T)$ .*

*Proof sketch.* Since  $\mathcal{W}$  is a standard measurable space we have that  $\mathcal{B}_{\mathcal{W}}$  is countably generated, i.e.  $\mathcal{B}_{\mathcal{W}} = \sigma(\mathcal{S})$  with a countable  $\mathcal{S} = \{D_n \mid n \in \mathbb{N}\} \subseteq \mathcal{B}_{\mathcal{W}}$ . If for example,  $\mathcal{W} = [0, 1]$ , which we can w.l.o.g. assume, then we could choose  $\mathcal{S} = \{[a, b] \mid a \leq b, a, b \in \mathbb{Q} \cap [0, 1]\}$ . We now define the following sequence of finite measurable partitions of  $\mathcal{W}$  inductively via:

$$\mathcal{E}_0 := \{\mathcal{W}\}, \quad \mathcal{E}_{n+1} := \left( \bigcup_{D \in \mathcal{E}_n} \{D \setminus D_n, D \cap D_n\} \right) \setminus \{\emptyset\}, \quad n \in \mathbb{N}.$$

We put  $\mathcal{B}_n := \sigma(\mathcal{E}_n)$ . Note that each  $\mathcal{E}_n$  is finite and for every  $n \in \mathbb{N}$ :

$$\mathcal{W} = \dot{\bigcup}_{D \in \mathcal{E}_n} D, \quad \mathcal{B}_n \subseteq \mathcal{B}_{n+1} \subseteq \mathcal{B}_{\mathcal{W}} = \sigma \left( \bigcup_{m \in \mathbb{N}} \mathcal{E}_m \right).$$

For  $D \in \mathcal{B}_{\mathcal{W}}$  we can define the map  $q_D : \mathcal{T} \rightarrow \mathbb{R}_{\geq 0}$  via:

$$q_D(t) := \frac{Q(W \in D | T = t)}{K(W \in D | T = t)} \cdot \mathbb{1}_{K(W \in D | T = t) > 0} = \begin{cases} \frac{Q(W \in D | T = t)}{K(W \in D | T = t)}, & \text{if } K(W \in D | T = t) > 0, \\ 0, & \text{if } K(W \in D | T = t) = 0. \end{cases}$$

Since  $Q(W \in D | T = t)$  and  $K(W \in D | T = t)$  are measurable in  $t$  for each fixed  $D$  we see that  $q_D$  is  $\mathcal{B}_{\mathcal{T}}\text{-}\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable. For  $n \in \mathbb{N}$  we now define:

$$G_n(w, t) := \sum_{D \in \mathcal{E}_n} \mathbb{1}_D(w) \cdot q_D(t),$$

and:

$$G(w, t) := \liminf_{n \in \mathbb{N}} G_n(w, t), \quad g(w|t) := G(w, t) \cdot \mathbb{1}_{G(w, t) < \infty}.$$

We immediately see that every  $G_n$  is a  $(\mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}})\text{-}\mathcal{B}_{\mathbb{R}_{\geq 0}}$ -measurable map. As a countable limit of measurable functions also  $G$  and  $g$  are  $\mathcal{B}_{\mathcal{W}} \otimes \mathcal{B}_{\mathcal{T}}$ -measurable. We claim that  $g$  is a Doob-Radon-Nikodym derivative of  $Q(W|T)$  w.r.t.  $K(W|T)$ . Since we already showed that  $g$  is jointly measurable we are left to show that for every  $t \in \mathcal{T}$  and  $D \in \mathcal{B}_{\mathcal{W}}$  we have:

$$Q(W \in D | T = t) = \int \mathbb{1}_D(w) \cdot g(w|t) K(W \in dw | T = t).$$

So in the following we can fix  $t \in \mathcal{T}$  and only indicate the dependence on  $t$  with an index:

$$G_n^t(w) := G_n(w, t), \quad G^t(w) := G(w, t).$$

Notice that  $G_n^t$  is  $\mathcal{B}_n$ -measurable for  $n \in \mathbb{N}$ . In the following we will use that by construction of the  $\mathcal{E}_n$  for  $D \in \mathcal{E}_n$  and  $m \geq n$  we have the disjoint union decompositions:

$$D = \dot{\bigcup}_{\substack{A \in \mathcal{E}_m \\ A \subseteq D}} A, \quad \mathcal{W} = \dot{\bigcup}_{D \in \mathcal{E}_n} \left( \dot{\bigcup}_{\substack{A \in \mathcal{E}_m \\ A \subseteq D}} A \right).$$

Let  $m \geq n$  then we get:

$$\begin{aligned}
& G_n^t(w) \\
&= \sum_{D \in \mathcal{E}_n} \left[ \frac{Q(W \in D|T=t)}{K(W \in D|T=t)} \cdot \mathbb{1}_{K(W \in D|T=t) > 0} \right] \cdot \mathbb{1}_D(w) \\
&= \sum_{D \in \mathcal{E}_n} \left[ \sum_{\substack{A \in \mathcal{E}_m \\ A \subseteq D}} \frac{Q(W \in A|T=t)}{K(W \in D|T=t)} \cdot \mathbb{1}_{K(W \in D|T=t) > 0} \right] \cdot \mathbb{1}_D(w) \\
&= \sum_{D \in \mathcal{E}_n} \left[ \sum_{\substack{A \in \mathcal{E}_m \\ A \subseteq D}} \frac{Q(W \in A|T=t)}{K(W \in D|T=t)} \cdot \mathbb{1}_{Q(W \in A|T=t) > 0} \cdot \mathbb{1}_{K(W \in D|T=t) > 0} \right] \cdot \mathbb{1}_D(w) \\
&\stackrel{Q \ll K}{=} \sum_{D \in \mathcal{E}_n} \left[ \sum_{\substack{A \in \mathcal{E}_m \\ A \subseteq D}} \frac{Q(W \in A|T=t)}{K(W \in D|T=t)} \cdot \mathbb{1}_{K(W \in A|T=t) > 0} \cdot \mathbb{1}_{K(W \in D|T=t) > 0} \right] \cdot \mathbb{1}_D(w) \\
&= \sum_{D \in \mathcal{E}_n} \left[ \sum_{\substack{A \in \mathcal{E}_m \\ A \subseteq D}} \frac{Q(W \in A|T=t)}{K(W \in A|T=t)} \cdot \frac{K(W \in A|T=t)}{K(W \in D|T=t)} \cdot \mathbb{1}_{K(W \in A|T=t) > 0} \cdot \mathbb{1}_{K(W \in D|T=t) > 0} \right] \cdot \mathbb{1}_D(w) \\
&= \sum_{A \in \mathcal{E}_m} \sum_{\substack{D \in \mathcal{E}_n \\ D \supseteq A}} \frac{Q(W \in A|T=t)}{K(W \in A|T=t)} \cdot \frac{K(W \in A|T=t)}{K(W \in D|T=t)} \cdot \mathbb{1}_{K(W \in A|T=t) > 0} \cdot \mathbb{1}_{K(W \in D|T=t) > 0} \cdot \mathbb{1}_D(w) \\
&= \sum_{A \in \mathcal{E}_m} \frac{Q(W \in A|T=t)}{K(W \in A|T=t)} \cdot \mathbb{1}_{K(W \in A|T=t) > 0} \underbrace{\left[ \sum_{\substack{D \in \mathcal{E}_n \\ D \supseteq A}} \frac{K(W \in A|T=t)}{K(W \in D|T=t)} \cdot \mathbb{1}_{K(W \in D|T=t) > 0} \cdot \mathbb{1}_D(w) \right]}_{= \mathbb{E}_t[\mathbb{1}_A | \mathcal{B}_n](w)} \\
&= \sum_{A \in \mathcal{E}_m} \left[ \frac{Q(W \in A|T=t)}{K(W \in A|T=t)} \cdot \mathbb{1}_{K(W \in A|T=t) > 0} \right] \cdot \mathbb{E}_t[\mathbb{1}_A | \mathcal{B}_n](w) \\
&= \mathbb{E}_t \left[ \sum_{A \in \mathcal{E}_m} \frac{Q(W \in A|T=t)}{K(W \in A|T=t)} \cdot \mathbb{1}_{K(W \in A|T=t) > 0} \cdot \mathbb{1}_A \middle| \mathcal{B}_n \right] (w) \\
&= \mathbb{E}_t [G_m^t | \mathcal{B}_n] (w).
\end{aligned}$$

Note that we use  $\mathbb{E}_t[\cdot | \mathcal{B}_n]$  to indicate conditional expectations w.r.t.  $K(W|T=t)$  and  $\mathcal{B}_n$ . For the first conditional expectation see [Kle20] Lem. 8.10. So we get that  $G_n^t$  is a version of  $\mathbb{E}_t[G_m^t | \mathcal{B}_n]$  for all  $m \geq n$ . This shows that  $(G_n^t)_{n \in \mathbb{N}}$  is a martingale attached to the filtration  $(\mathcal{B}_n)_{n \in \mathbb{N}}$  w.r.t.  $K(W|T=t)$ . Furthermore, we can show that  $(G_n^t)_{n \in \mathbb{N}}$  is uniformly integrable w.r.t.  $K(W|T=t)$ , see [Kle20] Ex. 7.39. By the convergence theorem for uniformly integrable martingales, see [Kle20] Thm. 11.7, we get that  $G_n^t$  also converges in  $L^1$  to  $G^t$  w.r.t.  $K(W|T=t)$  and that  $G_n^t$  is a version of  $\mathbb{E}_t[G^t | \mathcal{B}_n]$  for

all  $n \in \mathbb{N}$ . So for  $D \in \mathcal{E}_n$  the function  $\mathbb{1}_D \cdot G_n^t$  is a version of  $\mathbb{E}_t[\mathbb{1}_D \cdot G^t | \mathcal{B}_n]$ . Taking expectation values shows:

$$\mathbb{E}_t [\mathbb{1}_D \cdot G^t] = \mathbb{E}_t [\mathbb{E}_t[\mathbb{1}_D \cdot G^t | \mathcal{B}_n]] = \mathbb{E}_t [\mathbb{1}_D \cdot G_n^t] = Q(W \in D | T = t).$$

Since this holds for all  $D \in \mathcal{E}_n$  and all  $n \in \mathbb{N}$  it also holds for all  $D \in \mathcal{B}_\mathcal{W} = \sigma(\bigcup_{n \in \mathbb{N}} \mathcal{E}_n)$  and we get:

$$\int \mathbb{1}_D(w) \cdot G(w, t) K(W \in dw | T = t) = \mathbb{E}_t [\mathbb{1}_D \cdot G^t] = Q(W \in D | T = t).$$

Since  $Q(W | T = t)$  is a finite measure the set  $\{w \in \mathcal{W} | G(w, t) = \infty\}$  is a  $K(W | T = t)$ -null set and we can replace  $G$  by  $g$  under the integral. This shows the claim.  $\square$

### 2.3.4. Transition Probability Spaces

**Definition 2.3.23** (Transition probability space). *Consider measurable spaces  $\mathcal{T}$  and  $\mathcal{W}$  and a Markov kernel:*

$$K(W | T) : \mathcal{T} \dashrightarrow \mathcal{W}, \quad (D, t) \mapsto K(W \in D | T = t).$$

*Then we call the tuple  $(\mathcal{W} \times \mathcal{T}, K(W | T))$  a transition probability space. It generalizes the notion of probability space, which can be recovered by taking  $\mathcal{T} = *$ .*

**Definition 2.3.24** (Conditional random variables). *A measurable map:*

$$X : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{X}$$

*starting from a transition probability space  $(\mathcal{W} \times \mathcal{T}, K(W | T))$  is called conditional random variable. It generalizes the notion of random variables and can be considered a family of random variables (measurably) parameterized by  $t \in \mathcal{T}$ . For  $t \in \mathcal{T}$  we also define the measurable map:*

$$X_t : \mathcal{W} \rightarrow \mathcal{X}, \quad w \mapsto X_t(w) := X(w, t),$$

*which can be considered a random variable on the probability space  $(\mathcal{W}, K(W | T = t))$ .*

**Example 2.3.25** (Special conditional random variables of importance). *Let  $(\mathcal{W} \times \mathcal{T}, K(W | T))$  be a transition probability space. Then we denote by:*

1.  $T$  the canonical projection onto  $\mathcal{T}$ :

$$T := \text{pr}_\mathcal{T} : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{T}, \quad (w, t) \mapsto T(w, t) := t.$$

2.  $*$  the constant conditional random variable:

$$* : \mathcal{W} \times \mathcal{T} \rightarrow *, \quad (w, t) \mapsto *,$$

*where  $* := \{*\}$  is the one-point space.*

## 2.4. Constructing Markov Kernels from Others

### 2.4.1. Marginal Markov Kernels

**Definition 2.4.1** (Marginalizing Markov kernels). *Let*

$$K(X, Y|T) : \mathcal{T} \dashrightarrow \mathcal{X} \times \mathcal{Y}$$

be a Markov kernel in two variables. We can then define the marginal Markov kernels as follows:

$$K(X|T) : \mathcal{T} \dashrightarrow \mathcal{X}, \quad (A, t) \mapsto K(X \in A, Y \in \mathcal{Y}|T = t),$$

and:

$$K(Y|T) : \mathcal{T} \dashrightarrow \mathcal{Y}, \quad (B, t) \mapsto K(X \in \mathcal{X}, Y \in B|T = t).$$

**Example 2.4.2** (Marginal Markov kernels of discrete Markov kernels). *Let*

$$K(X, Y|T) : \mathcal{T} \dashrightarrow \mathcal{X} \times \mathcal{Y}$$

be a Markov kernel in two variables on discrete spaces and  $k_{X,Y|T}$  its mass function. We can then compute the marginal Markov kernels as follows:

$$k_{X|T}(x|t) = \sum_{y \in \mathcal{Y}} k_{X,Y|T}(x, y|t),$$

and:

$$k_{Y|T}(y|t) = \sum_{x \in \mathcal{X}} k_{X,Y|T}(x, y|t).$$

Note, by abuse of notation, for simplicity, we often omit the indices and write  $k(x|t)$  and  $k(y|t)$  instead and distinguish these two functions just by the use of the argument symbols  $x$  and  $y$ .

### 2.4.2. Product of Markov Kernels

**Definition 2.4.3** (Product of Markov kernels). *Consider two Markov kernels:*

$$Q(Z|Y, W, T) : \mathcal{Y} \times \mathcal{W} \times \mathcal{T} \dashrightarrow \mathcal{Z}, \quad K(W, U|T, X) : \mathcal{T} \times \mathcal{X} \dashrightarrow \mathcal{W} \times \mathcal{U}.$$

Then we define the product Markov kernel:

$$Q(Z|Y, W, T) \otimes K(W, U|T, X) : \mathcal{Y} \times \mathcal{T} \times \mathcal{X} \dashrightarrow \mathcal{Z} \times \mathcal{W} \times \mathcal{U},$$

using measurable sets  $E \subseteq \mathcal{Z} \times \mathcal{W} \times \mathcal{U}$  via:  $(E, (y, t, x)) \mapsto$

$$\int \int \mathbb{1}_E(z, w, u) Q(Z \in dz|Y = y, W = w, T = t) K((W, U) \in d(w, u)|T = t, X = x),$$

where the inner integration is over  $z \in \mathcal{Z}$  and the outer integration over  $(w, u) \in \mathcal{W} \times \mathcal{U}$ .<sup>5</sup>

<sup>5</sup>The integration ordering actually does not matter, which follows from Fubini's theorem, Theorem 2.4.7.

**Remark 2.4.4.** Note that in the notation of the product of Markov kernels we use the suggestive symbols, e.g.  $W$  in  $K(W, U|T, X)$  and  $W$  in  $Q(Z|Y, W, T)$  to indicate which variables will be “coupled” in the product. A more precise notation could indicate this, e.g. by indices on the product symbol like  $\otimes_{(W_1, W_2)}$  or similar. However, our shorthand notation should not lead to much ambiguity during this course. The rule of thumb is that the output of a kernel at the r.h.s. of the product is coupled to a matching input of the kernel on the l.h.s. of the product.

**Example 2.4.5** (Product of discrete Markov kernels). Let  $Q(Z|Y, W, T)$  and  $K(W|T, X)$  be two Markov kernels on finite spaces. Let  $P(Z, W|Y, T, X) := Q(Z|Y, W, T) \otimes K(W|T, X)$  be the product of Markov kernels and  $p, q, k$  the corresponding mass functions. Then we have:

$$p(z_i, w_k | y_s, x_l, t_j) = q(z_i | y_s, w_k, t_j) \cdot k(w_k | t_j, x_l),$$

which is just the product of mass functions. For the corresponding stochastic tensors  $\tilde{P}, \tilde{Q}, \tilde{K}$  we get that:

$$\tilde{P} = \tilde{Q} \odot_{W, T} \tilde{K}$$

is the entry-wise product/Hadamard product of tensors (reflecting the above formula, i.e. indices for  $w_k, t_j$  are the same in  $q$  and  $k$ ).

**Exercise 2.4.6.** Show that the product of Markov kernels is associative. Under which conditions can we commute Markov kernels in products? For this you can use Fubini’s theorem. See also the comments below.

**Theorem 2.4.7** (Fubini’s Theorem, [Kle20] Thm. 14.19). Let  $(\mathcal{X}, \mu)$  and  $(\mathcal{Y}, \nu)$  be two ( $\sigma$ -)finite measure spaces and  $f : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty]$  a  $(\mathcal{B}_{\mathcal{X}} \otimes \mathcal{B}_{\mathcal{Y}})$ - $\mathcal{B}_{[0, \infty]}$ -measurable map. Then we have the equalities:

$$\int \left( \int f(x, y) \mu(dx) \right) \nu(dy) = \int \left( \int f(x, y) \nu(dy) \right) \mu(dx) = \int f(x, y) d(\mu \otimes \nu)(x, y).$$

**Remark 2.4.8** (Conventions about integration order). For  $D \in \mathcal{B}_{\mathcal{X}} \otimes \mathcal{B}_{\mathcal{Y}}$  Fubini’s theorem says:

$$\int \left( \int \mathbb{1}_D(x, y) \mu(dx) \right) \nu(dy) = \int \left( \int \mathbb{1}_D(x, y) \nu(dy) \right) \mu(dx) = \int \mathbb{1}_D(x, y) d(\mu \otimes \nu)(x, y).$$

The integral notation hides the fact that  $\mu \otimes \nu$  and  $\nu \otimes \mu$  can only be identified as measures if we also swap the order of the spaces  $\mathcal{X}$  and  $\mathcal{Y}$ .  $\mu \otimes \nu$  lives on  $\mathcal{X} \times \mathcal{Y}$  and  $\nu \otimes \mu$  lives on  $\mathcal{Y} \times \mathcal{X}$ . In more precise terms, we would have:

$$(\mu \otimes \nu)(D) = (\nu \otimes \mu)(D^s),$$

where  $D^s := \{(y, x) \in \mathcal{Y} \times \mathcal{X} \mid (x, y) \in D\} \in \mathcal{B}_{\mathcal{Y}} \otimes \mathcal{B}_{\mathcal{X}}$ . That said, we will always make this swap implicitly, and just write:

$$Q(Z|Y, T) \otimes K(U|T, X) = K(U|T, X) \otimes Q(Z|Y, T),$$

if there is no variable with the same symbol that occurs on the left of the conditioning bar in one Markov kernel and on the right of the conditioning bar in the other Markov kernel. This is justified by Fubini's theorem and our implicit swap convention. This will not lead to much ambiguity, when interpreted as measures under the integral, as one can match the variables  $Z, U$ , etc., to their corresponding arguments  $z, u$ , etc., in our suggestive notations. In a similar sense we also identify:

$$\begin{aligned} K(W|T = t, X = x) &= K(W|X = x, T = t), \\ Q(X \in A, Y \in B|T) &= Q(Y \in B, X \in A|T). \end{aligned}$$

### 2.4.3. Composition of Markov Kernels

**Definition 2.4.9** (Composition of Markov kernels). *Consider two Markov kernels:*

$$Q(Z|Y, W, T) : \mathcal{Y} \times \mathcal{W} \times \mathcal{T} \dashrightarrow \mathcal{Z}, \quad K(W, U|T, X) : \mathcal{T} \times \mathcal{X} \dashrightarrow \mathcal{W} \times \mathcal{U}.$$

Then we define their composition:

$$Q(Z|Y, W, T) \circ K(W, U|T, X) : \mathcal{Y} \times \mathcal{T} \times \mathcal{X} \dashrightarrow \mathcal{Z},$$

using measurable sets  $C \subseteq \mathcal{Z}$  via:

$$(C, (y, t, x)) \mapsto \int Q(Z \in C|Y = y, W = w, T = t) K(W \in dw|T = t, X = x).$$

Note that we implicitly marginalized  $U$  out, i.e. in the composition we integrate over all variables (here:  $W$  and  $U$ ) from the right hand Markov kernel. As a notation we will also write:

$$\begin{aligned} &Q(Z \in C|Y = y, W, T = t) \circ K(W, U|T = t, X = x) \\ &:= (Q(Z|Y, W, T) \circ K(W, U|T, X))(C|(y, t, x)). \end{aligned}$$

**Remark 2.4.10.** *It is clear from the definitions 2.4.9, 2.4.3 and 2.4.1 that the composition:*

$$Q(Z|Y, W, T) \circ K(W, U|T, X)$$

is the  $Z$ -marginal of the product:

$$Q(Z|Y, W, T) \otimes K(W, U|T, X).$$

Furthermore, while the operation  $\otimes$  leaves all the variables of the second Markov kernel, here  $W$  and  $U$ , "intact", the operation  $\circ$  marginalizes them all out. One could also think of an intermediate operation that specifies which variables are marginalized out and which stays, e.g. using a symbol  $\circ_{(W,W)}$  to indicate marginalization of input  $W$  (of  $Q(Z|Y, W, T)$ ) over output  $W$  (from  $K(W, U|T, X)$ ). We will not further investigate this and will only use  $\otimes$  and  $\circ$  as described.



**Remark 2.4.11** (Composition of deterministic Markov kernels). *Consider measurable maps:*

$$X : \mathcal{T} \rightarrow \mathcal{X}, \quad Z : \mathcal{X} \rightarrow \mathcal{Z},$$

*and their composition  $Z \circ X$ . Then the composition of the corresponding Markov kernels satisfies:*

$$\delta(Z \circ X|T) = \delta(Z|X) \circ \delta(X|T),$$

*where  $\delta(Z \in C|X = x) := \mathbb{1}_C(Z(x))$  and  $\delta(X \in A|T = t) := \mathbb{1}_A(X(t))$ .*

*So the composition of Markov kernels extends the composition of functions.*

*Proof.*

$$\begin{aligned} \delta(Z \in C|X) \circ \delta(X|T = t) &= (\delta(Z|X) \circ \delta(X|T))(C|t) \\ &= \int \delta(Z \in C|X = x) \delta(X \in dx|T = t) \\ &= \int \mathbb{1}_{Z^{-1}(C)}(x) \delta(X \in dx|T = t) \\ &= \delta(X \in Z^{-1}(C)|T = t) \\ &= \mathbb{1}_{X^{-1}(Z^{-1}(C))}(t) \\ &= \mathbb{1}_C(Z(X(t))) \\ &= \delta(Z(X) \in C|T = t) \\ &= \delta(Z \circ X \in C|T = t) \\ &= \delta(Z \circ X|T)(C|t). \end{aligned}$$

□

**Example 2.4.12** (Composition of discrete Markov kernels). *Assume that all the spaces in definition 2.4.9 are discrete/finite and let  $P(Z|T) := Q(Z|W) \circ K(W|T)$  be the composition of Markov kernels. Let  $p, q, k$  denote the corresponding mass functions. Then we get:*

$$p(z_i|t_j) = \sum_k q(z_i|w_k) \cdot k(w_k|t_j).$$

*If  $\tilde{P}, \tilde{Q}, \tilde{K}$  are the corresponding stochastic matrices then we have that:*

$$\tilde{P} = \tilde{Q} \tilde{K},$$

*is just the usual matrix product. So in this case the composition of Markov kernels corresponds to matrix multiplication.*

#### 2.4.4. Push-Forward of Markov Kernels

**Definition 2.4.13** (Push-forward Markov kernel). *Let  $(\mathcal{W} \times \mathcal{T}, K(W|T))$  be a transition probability space and:*

$$X : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{X}$$

be a conditional random variable. Then we define the push-forward Markov kernel  $K(X|T)$  of  $K(W|T)$  w.r.t.  $X$  with symbols:

$$K(X|T) =: X_*K(W|T) =: K(X(W, T)|T),$$

via:

$$K(X|T) : \mathcal{T} \dashrightarrow \mathcal{X}, \quad (A, t) \mapsto K(X \in A|T = t) := K(W \in X_t^{-1}(A)|T = t),$$

where, again:

$$X_t^{-1}(A) = X^{-1}(A)_t := \{w \in \mathcal{W} \mid X(w, t) \in A\}.$$

**Remark 2.4.14.** We can also write push-forwards as compositions:

$$K(X|T) = \delta(X|W, T) \circ K(W|T),$$

where we define:

$$\delta(X \in A|W = w, T = t) := \mathbb{1}_A(X(w, t)) = \mathbb{1}_{X^{-1}(A)}(w, t).$$

**Remark 2.4.15.** For any Markov kernel

$$K(W|T) : \mathcal{T} \dashrightarrow \mathcal{W}$$

one can always extend it to include  $T = \text{pr}_{\mathcal{T}}$ :

$$K(W, T|T) : \mathcal{T} \dashrightarrow \mathcal{W} \times \mathcal{T}, \quad (E, t) \mapsto K((W, T) \in E|T = t) = K(W \in E_t|T = t),$$

where  $E_t = \{w \in \mathcal{W} \mid (w, t) \in E\}$ . Using Definition 2.4.3, we can also write this as:

$$K(W, T|T) = K(W|T) \otimes \delta(T|T),$$

where  $\delta(T \in D|T = t) := \mathbb{1}_D(t)$  for measurable  $D \subseteq \mathcal{T}$  and  $t \in \mathcal{T}$ .

### 2.4.5. Conditional Markov Kernels

**Definition/Theorem 2.4.16** (Disintegration of Markov kernels). Let  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $\mathcal{Z}$  be measurable spaces where  $\mathcal{X}$  and  $\mathcal{Y}$  are standard measurable spaces. Let

$$K(X, Y|Z) : \mathcal{Z} \dashrightarrow \mathcal{X} \times \mathcal{Y}$$

be a Markov kernel and  $K(Y|Z)$  its marginal Markov kernel given by  $K(Y \in B|Z) = K(X \in \mathcal{X}, Y \in B|Z)$ . Then there exists a Markov kernel (called conditional Markov kernel):

$$K(X|Y, Z) : \mathcal{Y} \times \mathcal{Z} \dashrightarrow \mathcal{X}$$

such that:

$$K(X, Y|Z) = K(X|Y, Z) \otimes K(Y|Z).$$

Furthermore,  $K(X|Y, Z)$  is essentially unique in the following sense: If  $Q(X|Y, Z)$  is another Markov kernel then we have:

$$K(X, Y|Z) = Q(X|Y, Z) \otimes K(Y|Z),$$

if and only if the measurable subset  $N$  of  $\mathcal{Y} \times \mathcal{Z}$  defined via:

$$N := \{(y, z) \in \mathcal{Y} \times \mathcal{Z} \mid \exists A \in \mathcal{B}_X : Q(X \in A|Y = y, Z = z) \neq K(X \in A|Y = y, Z = z)\}$$

is a  $K(Y|Z)$ -null set in  $\mathcal{Y} \times \mathcal{Z}$ .

**Remark 2.4.17.** If one further assumes certain continuity conditions for the conditional Markov kernel  $K(X|Y, Z)$  and that the marginal  $K(Y|Z)$  is strictly positive then the conditional Markov kernel can be fully identified, not just up to such  $K(Y|Z)$ -null sets. This is formalized in Lemma 2.4.23.

**Example 2.4.18** (Conditional Markov kernel for discrete Markov kernels). Consider a Markov kernel  $K(X, Y|Z)$  where all spaces are discrete and let  $k$  be the corresponding mass function. Then the marginal mass functions are given by:

$$k(y|z) = \sum_{x \in \mathcal{X}} k(x, y|z), \quad k(x|z) = \sum_{y \in \mathcal{Y}} k(x, y|z).$$

A conditional Markov kernel conditioned on  $Y$  can then be defined via the mass function:

$$k(x|y, z) := \begin{cases} \frac{k(x, y|z)}{k(y|z)} & \text{if } k(y|z) > 0, \\ k(x|z) & \text{if } k(y|z) = 0. \end{cases}^6$$

With this setting we then have for all (!) values  $x, y, z$ :

$$k(x, y|z) = k(x|y, z) \cdot k(y|z).$$

**Corollary 2.4.19** (Conditional probability distributions). Let  $X$  and  $Y$  be random variables on domain  $(\mathcal{W}, P(W))$  with standard measurable spaces  $\mathcal{X}, \mathcal{Y}$ , resp., as codomains. Then there always exist conditional probability distributions  $P(X|Y)$  and  $P(Y|X)$  that are Markov kernels satisfying:<sup>7</sup>

$$P(X, Y) = P(X|Y) \otimes P(Y), \quad P(X, Y) = P(Y|X) \otimes P(X).$$

Furthermore, these conditional probability distributions are essentially unique.

<sup>6</sup>Any value assignment for this spot is somewhat arbitrary as it almost surely does not occur. Typically this entry is defined to be 0. This is convenient but also problematic, as this would not normalize when summing over  $x \in \mathcal{X}$ . A proper alternative is to set it to be  $k(x|z)$  in this case.

<sup>7</sup>In the literature a conditional probability distribution that is also a Markov kernel would be called a *regular* version of a conditional probability distribution. Since in this lecture we will not encounter other versions we will just call this version here *conditional probability distribution*.

**Proofs - Disintegration of Markov Kernels** In this subsection we will give a proof for the existence and essential uniqueness of conditional Markov kernels. Another source for similar results can be found in [Kal17].

**Remark 2.4.20** (Existence of conditional Markov kernels). *If  $K(X, Y|Z)$  is a Markov kernel then we want  $K(X|Y, Z)$  such that:*

$$K(X, Y|Z) = K(X|Y, Z) \otimes K(Y|Z)$$

*holds. The heuristic here is to make use of Doob-Radon-Nikodym derivatives, see Theorem 2.3.17, for each  $A \in \mathcal{B}_X$ :*

$$K(X \in A|Y = y, Z = z) = \frac{K(X \in A, Y \in dy|Z = z)}{K(Y \in dy|Z = z)}(y).$$

*The problem is that they are only unique up to  $K(Y|Z)$ -null sets and might not be coordinated in such a way that  $K(X \in A|Y = y, Z = z)$  becomes a probability measure in  $A$  for every  $(y, z)$ . To ensure this we will take extra steps: We will first take the Doob-Radon-Nikodym derivative  $K(X \leq x|Y = y, Z = z)$  for rational points  $x \in \mathbb{Q}$  and then for general  $x \in \mathbb{R}$  put:*

$$K(X \leq x|Y = y, Z = z) = \inf_{m \in \mathbb{N}} K(X \leq \lceil x \rceil_m | Y = y, Z = z),$$

*where  $\lceil x \rceil_m := \frac{\lfloor mx+1 \rfloor}{m} \in \mathbb{Q}$  for  $m \in \mathbb{N}$ . This approach will work for  $K(Y|Z)$ -almost-all  $(y, z)$ . On the remaining points  $(y, z)$  we can then make a somewhat arbitrary choice, e.g. we can put:*

$$K(X \leq x|Y = y, Z = z) := K(X \leq x|Z = z).$$

*This will turn  $K(X \leq x|Y = y, Z = z)$  into a valid cumulative distribution function in  $x$  for all  $(y, z)$ , which then corresponds to a proper probability measure. One then checks that this  $K(X|Y, Z)$  is a desired conditional Markov kernel.*

**Theorem 2.4.21** (Existence of conditional Markov kernels). *Let  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$  be measurable spaces where  $\mathcal{X}$  is a standard measurable space and  $\mathcal{B}_Y$  is countably generated (e.g.  $\mathcal{Y}$  is also a standard measurable space). Let*

$$K(X, Y|Z) : \mathcal{Z} \dashrightarrow \mathcal{X} \times \mathcal{Y},$$

*be a Markov kernel in two variables. Then a conditional Markov kernel conditioned on  $Y$  given  $Z$ :*

$$K(X|Y, Z) : \mathcal{Y} \times \mathcal{Z} \dashrightarrow \mathcal{X},$$

*exists.*

*Proof.* Since  $\mathcal{X}$  is standard we can without loss of generality assume that  $\mathcal{X} = [0, 1]$ . For fixed  $A \in \mathcal{B}_X$  we have a finite transition measure  $K(X \in A, Y|Z)$  from  $\mathcal{Z}$  to  $\mathcal{Y}$ , which is absolute continuous w.r.t. the marginal  $K(Y|Z)$ , because of the inequality:

$$0 \leq K(X \in A, Y \in B|Z = z) \leq K(X \in \mathcal{X}, Y \in B|Z = z) = K(Y \in B|Z = z).$$

Since also  $\mathcal{B}_Y$  is countably generated, by Doob-Radon-Nikodym, see Theorem 2.3.17, we get a Doob-Radon-Nikodym derivative, i.e. a (jointly) measurable map:

$$g_A : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}_{\geq 0},$$

such that for all  $z \in \mathcal{Z}$  and  $B \in \mathcal{B}_Y$ :

$$K(X \in A, Y \in B | Z = z) = \int \mathbb{1}_B(y) \cdot g_A(y|z) K(Y \in dy | Z = z).$$

For  $x \in \mathcal{X}$  we will define:

$$G(x|y, z) := g_{[0,x]}(y, z).$$

As a next step we want to modify  $G(x|y, z)$  such that it becomes a cumulative distribution function in  $x$ , i.e. it corresponds to a probability distribution on  $\mathcal{X}$ . For this define  $\mathcal{X}_{\mathbb{Q}} := \mathcal{X} \cap \mathbb{Q}$ , which is countable and dense in  $\mathcal{X}$ . First note that:

$$S := \{(y, z) \in \mathcal{Y} \times \mathcal{Z} \mid G(1|y, z) \neq 1\}$$

is a measurable  $K(Y|Z)$ -null set. Then, for every pair  $x_1 < x_2$  in  $\mathcal{X}_{\mathbb{Q}}$  consider:

$$E_{(x_1, x_2)} := \{(y, z) \mid G(x_1|y, z) > G(x_2|y, z)\} \in \mathcal{B}_Y \otimes \mathcal{B}_Z.$$

Since we have the equations:

$$\begin{aligned} & \int \mathbb{1}_{E_{(x_1, x_2), z}}(y) \cdot G(x_1|y, z) K(Y \in dy | Z = z) \\ &= K(X \leq x_1, Y \in E_{(x_1, x_2), z} | Z = z) \\ &\stackrel{x_1 < x_2}{\leq} K(X \leq x_2, Y \in E_{(x_1, x_2), z} | Z = z) \\ &= \int \mathbb{1}_{E_{(x_1, x_2), z}}(y) \cdot G(x_2|y, z) K(Y \in dy | Z = z) \\ &\stackrel{G(x_2|y, z) < G(x_1|y, z)}{\leq} \int \mathbb{1}_{E_{(x_1, x_2), z}}(y) \cdot G(x_1|y, z) K(Y \in dy | Z = z) \end{aligned}$$

we necessarily have  $K(Y \in E_{(x_1, x_2), z} | Z = z) = 0$  for every  $z \in \mathcal{Z}$ .

Then  $E := S \cup \bigcup_{x_1 < x_2 \in \mathcal{X}_{\mathbb{Q}}} E_{(x_1, x_2)}$  is also a  $K(Y|Z)$ -null set in  $\mathcal{B}_Y \otimes \mathcal{B}_Z$ .

Now for  $x \in \mathcal{X}_{\mathbb{Q}}$  we can define:

$$D_x := \{(y, z) \mid G(x|y, z) < \inf_{n \in \mathbb{N}} G(x + 1/n|y, z)\} \in \mathcal{B}_Y \otimes \mathcal{B}_Z.$$

By the dominated convergence theorem (see [Kle20] Cor. 6.26) we get:

$$\begin{aligned} & \int \mathbb{1}_{D_{x,z}}(y) \cdot G(x|y, z) K(Y \in dy | Z = z) \\ &\stackrel{D_x}{\leq} \int \mathbb{1}_{D_{x,z}}(y) \cdot \inf_{n \in \mathbb{N}} G(x + \frac{1}{n}|y, z) K(Y \in dy | Z = z) \\ &= \inf_{n \in \mathbb{N}} \int \mathbb{1}_{D_{x,z}}(y) \cdot G(x + \frac{1}{n}|y, z) K(Y \in dy | Z = z) \\ &= \inf_{n \in \mathbb{N}} K(X \leq x + \frac{1}{n}, Y \in D_{x,z} | Z = z) \\ &= K(X \leq x, Y \in D_{x,z} | Z = z) \\ &= \int \mathbb{1}_{D_{x,z}}(y) \cdot G(x|y, z) K(Y \in dy | Z = z). \end{aligned}$$

So equality must hold, which then implies that:

$$\int \mathbb{1}_{D_{x,z}}(y) \cdot \underbrace{\left( \inf_{n \in \mathbb{N}} G(x + \frac{1}{n}|y, z) - G(x|y, z) \right)}_{>0 \text{ for } y \in D_{x,z}} K(Y \in dy|Z = z) = 0.$$

This shows that  $K(Y \in D_{x,z}|Z = z) = 0$  for all  $z \in \mathcal{Z}$ . So  $D := E \cup \bigcup_{x \in \mathcal{X}_{\mathbb{Q}}} D_x$  is again a  $K(Y|Z)$ -null set in  $\mathcal{B}_{\mathcal{Y}} \otimes \mathcal{B}_{\mathcal{Z}}$ .

So far, we got that  $G$ , when restricted to  $\mathcal{X}_{\mathbb{Q}} \times D^c$ , is jointly measurable in  $(y, z)$  for fixed  $x$  and monotone non-decreasing and continuous from above in  $x$  for fixed  $(y, z)$  with  $G(1|y, z) = 1$ . We now aim to extend  $G$  to  $\mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$ .

For  $x \in \mathcal{X} = [0, 1]$  and  $m \in \mathbb{N}$  put  $\lceil x \rceil_m := \min(1, \lfloor mx+1 \rfloor / m)$ . Then  $\lceil x \rceil_m \in [0, 1] \cap \mathbb{Q} = \mathcal{X}_{\mathbb{Q}}$ . The map  $x \mapsto \lceil x \rceil_m$  is measurable and for  $x \in [0, 1)$  we have:

$$x < \lceil x \rceil_m \leq x + \frac{1}{m}.$$

So  $\lceil 1 \rceil_m = 1$  and  $\lceil x \rceil_m \in \mathcal{X}_{\mathbb{Q}}$  converges to  $x \in \mathcal{X}$ ,  $x \neq 1$ , from above for  $m \rightarrow \infty$ . We then define for all  $(x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$ :

$$F(x|y, z) := \inf_{m \in \mathbb{N}} \{G(\lceil x \rceil_m|y, z)\} \cdot \mathbb{1}_{D^c}(y, z) + K(X \leq x|Z = z) \cdot \mathbb{1}_D(y, z).$$

It is clear that  $F$  is again jointly measurable in  $(y, z)$  for fixed  $x$  and agrees with  $G$  on  $\mathcal{X}_{\mathbb{Q}} \times D^c$  by construction. As a monotone approximation from above it is clearly continuous from above, monotone non-decreasing and satisfies  $F(1|y, z) = 1$  for all  $(y, z)$ . So for fixed  $(y, z)$  now  $F(\cdot|y, z)$  corresponds to a probability distribution  $K(X|Y = y, Z = z)$  on  $\mathcal{B}_{\mathcal{X}}$ , uniquely given by the defining relations on sets  $[0, x]$ :

$$F(x|y, z) =: K(X \leq x|Y = y, Z = z),$$

for all  $x \in \mathcal{X}$ .

Now define  $\mathcal{D} \subseteq \mathcal{B}_{\mathcal{X}}$  as the set of all  $A \in \mathcal{B}_{\mathcal{X}}$  that satisfy:

1. the map  $(y, z) \mapsto K(X \in A|Y = y, Z = z)$  is  $(\mathcal{B}_{\mathcal{Y}} \otimes \mathcal{B}_{\mathcal{Z}})$ - $\mathcal{B}_{\mathbb{R}}$ -measurable, and:
2. for all  $z \in \mathcal{Z}$  and  $B \in \mathcal{B}_{\mathcal{Y}}$  the following equation holds:

$$K(X \in A, Y \in B|Z = z) = \int \mathbb{1}_B(y) \cdot K(X \in A|Y = y, Z = z) K(Y \in dy|Z = z).$$

Since  $K(X, Y \in B|Z = z)$  and  $K(X|Y = y, Z = z)$  are measures in  $X$  the system  $\mathcal{D}$  is closed under countable disjoint unions. One can also check that  $\mathcal{D}$  is closed under complements and contains  $\mathcal{X} = [0, 1]$ . So  $\mathcal{D}$  is a Dynkin system. We already know that for  $x \in \mathcal{X}_{\mathbb{Q}}$  the map  $(y, z) \mapsto K(X \leq x|Y = y, Z = z) = F(x|y, z)$  is measurable. Since for  $x \in \mathcal{X}_{\mathbb{Q}}$  and every  $B \in \mathcal{B}_{\mathcal{Y}}$ ,  $z \in \mathcal{Z}$ , we have:

$$\mathbb{1}_B(y) \cdot K(X \leq x|Y = y, Z = z) = \mathbb{1}_B(y) \cdot G(x|y, z)$$

up to the  $K(Y|Z = z)$ -null set  $D_z$  we already get for those  $x \in \mathcal{X}_{\mathbb{Q}}$ :

$$K(X \leq x, Y \in B|Z = z) = \int \mathbb{1}_B(y) \cdot K(X \leq x|Y = y, Z = z) K(Y \in dy|Z = z).$$

This shows that  $\mathcal{E} := \{[0, x] \mid x \in \mathcal{X}_{\mathbb{Q}}\} \subseteq \mathcal{D}$ . Since  $\mathcal{E}$  is closed under finite intersections Dynkin's lemma (see [Kle20] Thm. 1.19) implies:

$$\mathcal{B}_{\mathcal{X}} = \sigma(\mathcal{E}) \subseteq \mathcal{D}.$$

This shows that the two conditions hold for all  $A \in \mathcal{B}_{\mathcal{X}}$  and thus that  $K(X|Y, Z)$  is the desired conditional Markov kernel.  $\square$

**Lemma 2.4.22** (Essential uniqueness). *If we have Markov kernels:*

$$P(X|Y, Z), Q(X|Y, Z) : \mathcal{Y} \times \mathcal{Z} \dashrightarrow \mathcal{X},$$

and

$$K(Y|Z) : \mathcal{Z} \dashrightarrow \mathcal{Y}$$

with any measurable spaces  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$  such that:

$$P(X|Y, Z) \otimes K(Y|Z) = Q(X|Y, Z) \otimes K(Y|Z),$$

then for every  $A \in \mathcal{B}_{\mathcal{X}}$  the set:

$$N_A := \{(y, z) \in \mathcal{Y} \times \mathcal{Z} \mid P(X \in A|Y = y, Z = z) \neq Q(X \in A|Y = y, Z = z)\}$$

is a measurable  $K(Y|Z)$ -null set.

If, furthermore,  $\mathcal{X}$  is countably generated, e.g. a standard measurable space, then also  $N := \bigcup_{A \in \mathcal{B}_{\mathcal{X}}} N_A$  is a measurable  $K(Y|Z)$ -null set.

*Proof.* For fixed  $A \in \mathcal{B}_{\mathcal{X}}$  both  $P(X \in A|Y, Z)$  and  $Q(X \in A|Y, Z)$  can be considered a Doob-Radon-Nikodym derivative of the same transition measure  $M_A(Y|Z)$  given by:

$$\begin{aligned} M_A(Y \in B|Z = z) &:= \int \mathbb{1}_B(y) \cdot P(X \in A|Y, Z) K(Y \in dy|Z = z) \\ &= (P(X \in A|Y, Z) \otimes K(Y|Z))(B|z) \\ &= (Q(X \in A|Y, Z) \otimes K(Y|Z))(B|z) \\ &= \int \mathbb{1}_B(y) \cdot Q(X \in A|Y, Z) K(Y \in dy|Z = z). \end{aligned}$$

The uniqueness statement then follows from that of Doob-Radon-Nikodym derivatives, see Lemma 2.3.21. If now  $\mathcal{B}_{\mathcal{X}}$  is countably generated then  $\mathcal{B}_{\mathcal{X}} = \sigma(\mathcal{A})$  with a countable set  $\mathcal{A}$  that is closed under finite intersections, e.g.  $\mathcal{B}_{[0,1]} = \sigma(\{[0, c] \mid c \in [0, 1] \cap \mathbb{Q}\})$ . One then puts  $M := \bigcup_{A \in \mathcal{A}} N_A$ , which is, as countable union of  $K(Y|Z)$ -null sets, a  $K(Y|Z)$ -null set. Then one can define:

$$\mathcal{D} := \{A \in \mathcal{B}_{\mathcal{X}} \mid \forall (y, z) \in M^c : P(X \in A|Y = y, Z = z) = Q(X \in A|Y = y, Z = z)\}.$$

One easily sees that  $\mathcal{D}$  is closed under complements, countable disjoint unions and contains  $\mathcal{X}$ . This shows that  $\mathcal{D}$  is a Dynkin system (aka  $\lambda$ -system). Furthermore, we have:  $\mathcal{A} \subseteq \mathcal{D}$  and that  $\mathcal{A}$  is closed under finite intersections. By Dynkin's lemma we get that:

$$\mathcal{B}_{\mathcal{X}} = \sigma(\mathcal{A}) \subseteq \mathcal{D}.$$

This shows that  $N = \bigcup_{A \in \mathcal{B}_{\mathcal{X}}} N_A \subseteq M$ , thus  $N = M$  which is a measurable  $K(Y|Z)$ -null set.  $\square$

We now want to prove that we can recover from the ambiguity of the null sets for conditional Markov kernel under continuity assumptions and strictly positive marginals.

**Lemma 2.4.23** (Uniqueness for continuous conditional Markov kernels and strictly positive marginals). *Let  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $\mathcal{Z}$  be Polish spaces endowed with their Borel- $\sigma$ -algebra and:*

$$P(X|Y, Z), Q(X|Y, Z) : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{X}),$$

*two continuous Markov kernels, where  $\mathcal{P}(\mathcal{X})$  carries any Hausdorff topology  $\mathcal{T}_{\mathcal{X}} \subseteq \mathcal{B}_{\mathcal{P}(\mathcal{X})}$ , e.g. the weak\*-topology. Let*

$$K(Y|Z) : \mathcal{Z} \dashrightarrow \mathcal{Y}$$

*be a Markov kernel that is strictly positive (on non-empty open subsets of  $\mathcal{Y}$ ). If we have the equality of Markov kernels  $\mathcal{Z} \dashrightarrow \mathcal{X} \times \mathcal{Y}$ :*

$$P(X|Y, Z) \otimes K(Y|Z) = Q(X|Y, Z) \otimes K(Y|Z),$$

*then we already have the equality of Markov kernels:*

$$P(X|Y, Z) = Q(X|Y, Z).$$

*Proof.* Consider the set:

$$\Delta_{\mathcal{P}(\mathcal{X})} := \{(P, P) \in \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X}) \mid P \in \mathcal{P}(\mathcal{X})\},$$

which is a closed subset of  $\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X})$  because  $\mathcal{P}(\mathcal{X})$  is Hausdorff. Then the set:

$$\begin{aligned} N &:= \{(y, z) \in \mathcal{Y} \times \mathcal{Z} \mid P(X|Y = y, Z = z) \neq Q(X|Y = y, Z = z)\} \\ &= \{(y, z) \in \mathcal{Y} \times \mathcal{Z} \mid (P(X|Y = y, Z = z), Q(X|Y = y, Z = z)) \notin \Delta_{\mathcal{P}(\mathcal{X})}\} \\ &= (P(X|Y, Z), Q(X|Y, Z))^{-1}(\Delta_{\mathcal{P}(\mathcal{X})}^c), \end{aligned}$$

is an open subset of  $\mathcal{Y} \times \mathcal{Z}$  as both Markov kernels are continuous. By the essential uniqueness from Lemma 2.4.22 we know that for all  $z \in \mathcal{Z}$  we have:

$$K(Y \in N^z | Z = z) = 0.$$

The fact that the section  $N^z$  is open in  $\mathcal{Y}$  and that  $K(Y|Z)$  is strictly positive implies that either  $K(Y \in N^z | Z = z) > 0$  or that  $N^z = \emptyset$ . Since the former was ruled out by the essential uniqueness we get  $N^z = \emptyset$  for all  $z \in \mathcal{Z}$  and thus  $N = \emptyset$ . This shows the claim:

$$P(X|Y, Z) = Q(X|Y, Z).$$

$\square$



## 2.5. Conditional Independence

### 2.5.1. Independence for Random Variables

**Motivation 2.5.1.** *If we throw two dice with outcome values  $X$  and  $Y$ , resp., then knowing the value of  $Y$  does not give us any information about the value of  $X$ , and vice versa. We say that  $X$  and  $Y$  are independent from each other. We will formalize this intuition for all random variables in the following.*

**Definition 2.5.2** (Independence of two random variables). *Let  $(\mathcal{W}, P(W))$  be a probability space and  $X : \mathcal{W} \rightarrow \mathcal{X}$  and  $Y : \mathcal{W} \rightarrow \mathcal{Y}$  be two random variables. We say that  $X$  and  $Y$  are independent if the following equation holds:*

$$P(X, Y) = P(X) \otimes P(Y),$$

where  $P(X, Y)$  is the joint and  $P(X)$  and  $P(Y)$  are the corresponding marginal distributions. In symbols we would write this as:

$$X \underset{P(W)}{\perp\!\!\!\perp} Y.$$

**Lemma 2.5.3.** *Let  $(\mathcal{W}, P(W))$  be a probability space,  $\mathcal{X}$  and  $\mathcal{Y}$  standard measurable spaces and  $X : \mathcal{W} \rightarrow \mathcal{X}$  and  $Y : \mathcal{W} \rightarrow \mathcal{Y}$  be two random variables. Then the following statements are equivalent:*

1.  $X \underset{P(W)}{\perp\!\!\!\perp} Y$ .
2.  $P(X, Y) = P(X) \otimes P(Y)$ .
3. *There exists a probability distribution  $Q(X)$  such that:*

$$P(X, Y) = Q(X) \otimes P(Y).$$

4.  $P(X|Y) = P(X)$  holds  $P(Y)$ -almost-surely, where  $P(X|Y)$  is a version of a conditional probability distribution from Cor. 2.4.19.
5. *For all  $A \in \mathcal{B}_{\mathcal{X}}$  we have:*

$$\mathbb{E}[\mathbb{1}_A(X)|Y] = \mathbb{E}[\mathbb{1}_A(X)] \quad P(W)\text{-a.s.}$$

6. *For all  $A \in \mathcal{B}_{\mathcal{X}}$  and  $B \in \mathcal{B}_{\mathcal{Y}}$  we have:*

$$P(X \in A, Y \in B) = P(X \in A) \cdot P(Y \in B).$$

*Proof.* Exercise. □

**Exercise 2.5.4.** *Reformulate the statements in Lemma 2.5.3 for the case we either have mass functions (discrete case) or densities w.r.t. a product measure, e.g. the Lebesgue measure (absolute continuous case).*

We can generalize the notion of independence to arbitrary families of random variables:

**Definition 2.5.5** (Mutual independence of families of random variables). *Let  $(\mathcal{W}, P(W))$  be a probability space and  $I$  an (index) set. For  $i \in I$  let  $X_i : \mathcal{W} \rightarrow \mathcal{X}_i$  be a random variable. We say that  $(X_i)_{i \in I}$  is (mutually/jointly) independent if for all two disjoint subsets  $J_1 \dot{\cup} J_2 \subseteq I$  we have the independence:*

$$(X_{j_1})_{j_1 \in J_1} \perp\!\!\!\perp_{P(W)} (X_{j_2})_{j_2 \in J_2}.$$

**Exercise 2.5.6** (Mutual independence for finite tuples of random variables). *A finite tuple of random variables  $(X_1, \dots, X_n)$  is mutually independent if and only if:*

$$P(X_1, \dots, X_n) = P(X_1) \otimes \dots \otimes P(X_n).$$

**Exercise 2.5.7.** *Let  $(\mathcal{W}, P(W))$  be a probability space and  $I$  an arbitrary index set. For  $i \in I$  let  $X_i : \mathcal{W} \rightarrow \mathcal{X}_i$  be a random variable. The following statements are equivalent:*

1.  $(X_i)_{i \in I}$  is (mutually/jointly) independent.
2. For every finite disjoint subsets  $J_1, J_2 \subseteq I$  we have the independence:

$$(X_{j_1})_{j_1 \in J_1} \perp\!\!\!\perp_{P(W)} (X_{j_2})_{j_2 \in J_2}.$$

3. For every finite subset  $J \subseteq I$  and  $i \in I \setminus J$  we have:

$$X_i \perp\!\!\!\perp_{P(W)} (X_j)_{j \in J}.$$

4. For every finite subset  $J \subseteq I$  we have:

$$P((X_j)_{j \in J}) = \bigotimes_{j \in J} P(X_j).$$

5. We have the equality:

$$P((X_i)_{i \in I}) = \bigotimes_{i \in I} P(X_i),$$

where  $\bigotimes_{i \in I} P(X_i)$  is the product measure on  $\mathcal{X} = \prod_{i \in I} \mathcal{X}_i$ , which is determined by the corresponding products on its finite marginals via the extension theorem of Ionescu-Tulcea, see [IT49, Lam87] and theorem A.10.2.

## 2.5.2. Conditional Independence for Random Variables

**Motivation 2.5.8.** 1. Consider two independent coin flips with outcome variables  $X$  and  $Y$ , resp., with values in  $\{0, 1\}$ , and  $Z := X + Y \in \{0, 1, 2\}$ . If the value of  $Z$  is known, say  $Z = 1$ , then revealing the value of  $Y$ , say  $Y = 0$ , provides us with all the information to fully determine the value of  $X$ , here  $X = 0$ . This is despite the fact that  $X$  and  $Y$  were assumed to be independent. This means that conditioning on a third variable  $Z$  can destroy independence. In this case, we say that  $X$  and  $Y$  are dependent conditioned on  $Z$ . Summarized in symbols we have:

$$X \underset{P(W)}{\perp\!\!\!\perp} Y, \quad \text{but:} \quad X \not\underset{P(W)}{\perp\!\!\!\perp} Y | Z.$$

2. Now consider three (mutually) independent coin flips  $X, W, U$  with values in  $\{0, 1\}$ . Let  $Z := X + W$  and  $Y := Z + U = X + W + U$ . If we knew the value of  $Y$ , say  $Y = 0$ , then we would have information about the values of  $X$  as well, here  $X = 0$ . This shows that  $X$  and  $Y$  can not be independent random variables. If, in contrast, we would first reveal the value of  $Z$ , say  $Z = 1$ , then the value of  $X$  might be restricted by the value of  $Z$ , but also revealing  $Y$  would not give us any additional information about the value of  $X$ . The reason is that  $Y = Z + U$  and  $U$  is independent of  $X, W$  and  $Z = X + W$ . So, even though  $X$  and  $Y$  are dependent, when conditioned on  $Z$  they become independent, as there is no additional information gained about each others value, when revealing the other. Summarized in symbols we have:

$$X \not\underset{P(W)}{\perp\!\!\!\perp} Y, \quad \text{but:} \quad X \underset{P(W)}{\perp\!\!\!\perp} Y | Z.$$

We now want to formalize conditional independence for random variables.

**Remark 2.5.9** (Conditional independence). In contrast to (unconditional) independence, see Definition 2.5.2, possible definitions of conditional independence come with many more subtleties, due to their interplay with conditional probability distributions or conditional expectations. Such definitions can in general be non-equivalent. However, if we restrict ourselves to standard measurable spaces the subtleties can be resolved and the definitions become equivalent. This is the reason that in the following we will only state conditional independence for standard measurable spaces. We will make a clearer choice later for conditional independence of conditional random variables.

**Definition/Lemma 2.5.10** (Conditional independence for random variables). Let  $(\mathcal{W}, P(W))$  be a probability space and  $\mathcal{X}, \mathcal{Y}$  and  $\mathcal{Z}$  standard measurable spaces, and  $X : \mathcal{W} \rightarrow \mathcal{X}$  and  $Y : \mathcal{W} \rightarrow \mathcal{Y}$  and  $Z : \mathcal{W} \rightarrow \mathcal{Z}$  be three random variables. We then say that  $X$  is independent of  $Y$  conditioned on  $Z$  if any of the following equivalent conditions holds:

1.  $P(X, Y | Z) = P(X | Z) \otimes P(Y | Z)$  holds  $P(Z)$ -a.s., where  $P(X, Y | Z), P(X | Z), P(Y | Z)$  are versions of conditional probability distributions from Cor. 2.4.19.

2.  $P(X|Y, Z) = P(X|Z)$  holds  $P(Y, Z)$ -a.s., where  $P(X|Y, Z)$ ,  $P(X|Z)$  from Cor. 2.4.19.

3. There exists a Markov kernel  $Q(X|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}$  such that:

$$P(X, Y, Z) = Q(X|Z) \otimes P(Y, Z).$$

4. For every  $A \in \mathcal{B}_X$  and  $B \in \mathcal{B}_Y$  we have:

$$\mathbb{E}[\mathbb{1}_A(X) \cdot \mathbb{1}_B(Y)|Z] = \mathbb{E}[\mathbb{1}_A(X)|Z] \cdot \mathbb{E}[\mathbb{1}_B(Y)|Z] \quad P(W)\text{-a.s.}$$

5. For every  $A \in \mathcal{B}_X$  we have:

$$\mathbb{E}[\mathbb{1}_A(X)|Y, Z] = \mathbb{E}[\mathbb{1}_A(X)|Z] \quad P(W)\text{-a.s.}$$

In those cases, in symbols we write:

$$X \underset{P(W)}{\perp\!\!\!\perp} Y | Z.$$

*Proof.* Exercise. □

**Exercise 2.5.11.** Restate all the statements in Definition/Lemma 2.5.10 for discrete random variables in terms of mass functions.

Conditional independence for random variables satisfies the following rules:

**Theorem 2.5.12** (Separoid axioms for conditional independence for random variables, see [Daw01]). Let  $(W, P(W))$  be a probability space and  $X, Y, Z$  and  $U$  random variables taking values in standard measurable spaces  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$  and  $\mathcal{U}$ , respectively. Then we have the following rules:

1. *Redundancy:* If  $U = \varphi(X)$  a.s. is a measurable function of  $X$ , e.g.  $U = X$ , then:

$$U \underset{P(W)}{\perp\!\!\!\perp} Y | X.$$

2. *Symmetry:*

$$X \underset{P(W)}{\perp\!\!\!\perp} Y | Z \quad \implies \quad Y \underset{P(W)}{\perp\!\!\!\perp} X | Z.$$

3. *Decomposition:*

$$X \underset{P(W)}{\perp\!\!\!\perp} Y, U | Z \quad \implies \quad X \underset{P(W)}{\perp\!\!\!\perp} U | Z.$$

4. *Weak Union:*

$$X \underset{P(W)}{\perp\!\!\!\perp} Y, U | Z \quad \implies \quad X \underset{P(W)}{\perp\!\!\!\perp} Y | U, Z.$$

5. *Contraction:*

$$\left( X \perp\!\!\!\perp_{P(W)} U \mid Z \right) \wedge \left( X \perp\!\!\!\perp_{P(W)} Y \mid U, Z \right) \implies X \perp\!\!\!\perp_{P(W)} Y, U \mid Z.$$

*Proof.* Exercise. □

**Remark 2.5.13.** *The separoid axioms, see Theorem 2.5.12, also hold true for random variables that map into general (non-standard) measurable spaces if one restricts oneself to the definition of conditional independence only involving conditional expectations (rather than conditional probabilities or Markov kernels) in Definition 2.5.10.*

**Exercise 2.5.14.** *Assume that the random variables  $X, Y, Z, U$  have a joint density  $p$  w.r.t. some product measure (or a joint mass function) such that  $p(y, u|z) > 0$  for all values  $y, u, z$ . Show that we then also have the following intersection rule:*

$$\left( X \perp\!\!\!\perp_{P(W)} U \mid Y, Z \right) \wedge \left( X \perp\!\!\!\perp_{P(W)} Y \mid U, Z \right) \implies X \perp\!\!\!\perp_{P(W)} Y, U \mid Z.$$

### 2.5.3. Conditional Independence for Conditional Random Variables

**Motivation 2.5.15.** *Assume that we are given a statistical model  $P(X|\Theta)$  and a statistic  $S = S(X)$ , which is a measurable function of  $X$ . Often one wants to find such an  $S$  such that the choice of parameter  $\Theta = \theta$  has no “influence” on the probability distribution of  $X$  when  $S$  is provided. Such a statistic is usually called a sufficient statistic of  $X$  w.r.t.  $P(X|\Theta)$ . In symbols we want  $S$  such that:*

$$X \perp\!\!\!\perp \Theta \mid S.$$

*However, the parameter variable  $\Theta$  here is not a proper random variable as we have no distribution  $P(\Theta)$  specified over it. Still such a conditional independence statement makes sense. We thus want to formalize a notion of conditional independence for conditional random variables. We follow the definition of [For21]. Other approaches can be found in [Daw79, Daw80, Daw01, CD17, RERS23, FM20].*

**Motivation 2.5.16.** *Consider a probabilistic program with input variables  $T, S$  and output variables  $X, Y, Z$ . Whenever the program is given  $T$  and  $S$  as input, it internally samples  $U, E \sim U[0, 1]$  uniformly and independently from a random number generator, then calculates:*

$$X := T + S + U, \quad Y := 5 \cdot S + E, \quad Z := X \cdot Y,$$

*and, finally, outputs  $X, Y$  and  $Z$ . Even though, the input  $T$  and  $S$  is provided by the user and is not considered a random variable, we can reason about the fact that “Output  $Y$  only depends on the input  $S$  and not on the input  $T$ .” We want to formalize such conditional independence mathematically in order to be able to write this as:*

$$Y \perp\!\!\!\perp T \mid S.$$

**Definition 2.5.17** (Conditional independence for conditional random variables). Let  $(\mathcal{W} \times \mathcal{T}, K(W|T))$  be a transition probability space with Markov kernel:

$$K(W|T) : \mathcal{T} \dashrightarrow \mathcal{W}.$$

Consider conditional random variables:

$$X : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{X}, \quad Y : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{Y}, \quad Z : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{Z}.$$

We say that  $X$  is independent of  $Y$  conditioned on  $Z$  w.r.t.  $K(W|T)$ , in symbols:

$$X \perp\!\!\!\perp_{K(W|T)} Y \mid Z,$$

if there exists a Markov kernel:

$$Q(X|Z) : \mathcal{Z} \dashrightarrow \mathcal{X},$$

such that:

$$K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T),$$

where  $K(Y, Z|T)$  is the marginal of  $K(X, Y, Z|T)$ .<sup>8</sup>

As a special case, we define:

$$X \perp\!\!\!\perp_{K(W|T)} Y \quad : \iff \quad X \perp\!\!\!\perp_{K(W|T)} Y \mid *.$$

**Notation 2.5.18** (Essential uniqueness). The Markov kernel  $Q(X|Z)$  appearing in the conditional independence  $X \perp\!\!\!\perp_{K(W|T)} Y \mid Z$  in definition 2.5.17 is then a version of a conditional Markov kernel  $K(X|Y, Z, T)$  and is thus essentially unique in the sense of 2.4.22. We will use the following suggestive notation for it:

$$K(X|\cancel{T}, Y, Z) := Q(X|Z),$$

or similarly with crossed variables in different order. So we have in case of  $X \perp\!\!\!\perp_{K(W|T)} Y \mid Z$ :

$$K(X, Y, Z|T) = K(X|\cancel{T}, Y, Z) \otimes K(Y, Z|T).$$

Note that  $K(X|\cancel{T}, Y, Z)$  is a version of the conditional Markov kernel  $K(X|Y, Z, T)$  and does not depend on arguments  $y$  and  $t$ .

**Remark 2.5.19** (Conditional independence includes conditional independence from  $T$ ). We have the equivalence:

$$X \perp\!\!\!\perp_{K(W|T)} Y \mid Z \quad \iff \quad X \perp\!\!\!\perp_{K(W|T)} T, Y \mid Z,$$

where  $T : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{T}$ ,  $(w, t) \mapsto t$ , is the canonical projection map.

---

<sup>8</sup>For the equation  $K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T)$  to hold it is sufficient to check that for all  $t \in \mathcal{T}$ ,  $A \in \mathcal{B}_X$ ,  $B \in \mathcal{B}_Y$  and  $C \in \mathcal{B}_Z$  we have:

$$K(X \in A, Y \in B, Z \in C|T = t) = \int_C \int_B Q(X \in A|Z = z) K(Y \in dy, Z \in dz|T = t).$$

*Proof.*

$$\begin{aligned}
& X \underset{K(W|T)}{\perp\!\!\!\perp} Y \mid Z \\
\iff & \exists Q(X|Z) : K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T) \\
\iff & \exists Q(X|Z) : K(X, T, Y, Z|T) = Q(X|Z) \otimes \underbrace{K(Y, Z|T) \otimes \delta(T|T)}_{K(T, Y, Z|T)} \\
\iff & X \underset{K(W|T)}{\perp\!\!\!\perp} T, Y \mid Z.
\end{aligned}$$

The middle implication “ $\implies$ ” follows by taking the product with  $\delta(T|T)$ , and the reverse implication “ $\impliedby$ ” by marginalizing out  $T$ , i.e. via  $\delta(T \in \mathcal{T}|T) = 1$ .  $\square$

**Remark 2.5.20** (How to find  $Q(X|Z)$  and check for conditional independence?). *In case we have the conditional independence:*

$$X \underset{K(W|T)}{\perp\!\!\!\perp} Y \mid Z,$$

*we then get by definition:*

$$K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T),$$

*for some Markov kernel  $Q(X|Z)$ . This implies for all  $t \in \mathcal{T}$  the equation:*

$$K(X, Z|T = t) = Q(X|Z) \otimes K(Z|T = t).$$

*This means that  $Q(X|Z)$  is a version of the conditional probability distribution  $K(X|Z, T = t)$  for all  $t \in \mathcal{T}$  at once, and, in addition, it is also functionally not dependent on  $t$ . So for fixed  $t_0 \in \mathcal{T}$  the conditional  $K(X|Z, T = t_0)$  can be changed on a  $K(Z|T = t_0)$ -null-set such that it agrees with  $Q(X|Z)$ . So it is reasonable to test out versions of  $K(X|Z, T = t_0)$  for  $Q(X|Z)$ . To summarize, we have the following equivalence between:*

1.  $X \underset{K(W|T)}{\perp\!\!\!\perp} Y \mid Z,$
2. *There exist  $t_0 \in \mathcal{T}$  and a (regular) version of the conditional probability distribution  $K(X|Z, T = t_0)$  such that for all  $t \in \mathcal{T}$ :*

$$K(X, Y, Z|T = t) = K(X|Z, T = t_0) \otimes K(Y, Z|T = t).$$

*Note that in the last expression the middle term has the fixed  $t_0$  and the outer two terms have varying  $t \in \mathcal{T}$ .*

*This equivalence allows us to narrow our search space for  $Q(X|Z)$  to such conditional probability distributions.*

**Example 2.5.21** (Conditional independence for discrete conditional random variables). *Let the situation be like in definition 2.5.17 and assume all spaces to be countable and discrete. Let  $k$  be the mass function for  $K(X, Y, Z|T)$ . Then we have:*

$$X \underset{K(W|T)}{\perp\!\!\!\perp} Y | Z,$$

*if and only if there is a probability mass function  $q$  such that for all values  $x, y, z, t$ :*

$$k(x, y, z|t) = q(x|z) \cdot k(y, z|t).$$

*Note that in this case  $q(x|z)$  is a version of  $k(x|z, t)$  for all  $t \in \mathcal{T}$  at once, but that is also independent of  $t$ . We can use this knowledge to find such a proposal  $q(x|z)$  as follows.*

*If there exists a  $t_0 \in \mathcal{T}$  such that  $k(z|t_0) > 0$  for all  $z \in \mathcal{Z}$  then the conditional  $k(x|z, t_0)$  is uniquely given and equal to  $\frac{k(x, z|t_0)}{k(z|t_0)}$ .  $k(x|z, t_0)$  would then necessarily agree with  $q(x|z)$  in case of the conditional independence. So, if there exists a  $t_0 \in \mathcal{T}$  such that  $k(z|t_0) > 0$  for all  $z \in \mathcal{Z}$  then we get the following equivalence:*

$$X \underset{K(W|T)}{\perp\!\!\!\perp} Y | Z \quad \iff \quad \forall x, y, z, t : \quad k(x, y, z|t) = k(x|z, t_0) \cdot k(y, z|t).$$

*Again, note that in the last expression the middle term has the fixed  $t_0$  and the outer two terms have varying  $t \in \mathcal{T}$ .*

This example can be generalized.

**Theorem 2.5.22** (Conditional independence for conditional random variables with density). *Let  $\mu_X, \mu_Y$  and  $\mu_Z$  reference measures on  $\mathcal{X}, \mathcal{Y}$  and  $\mathcal{Z}$ , resp., and  $\mu := \mu_X \otimes \mu_Y \otimes \mu_Z$  the product measure on  $\mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$ . Assume that  $K(X, Y, Z|T)$  has a Doob-Radon-Nikodym derivative  $k$  w.r.t.  $\mu$  and let  $t_0 \in \mathcal{T}$  be a fixed value. Then we have the implication:*

$$\forall t \in \mathcal{T} \forall_\mu x, y, z. \quad k(x, y, z|t) = k(x|z, t_0) \cdot k(y, z|t) \quad \implies \quad X \underset{K(W|T)}{\perp\!\!\!\perp} Y | Z,$$

*where  $\forall_\mu$  means “for  $\mu$ -almost-all”. If for  $\mu_Z$ -almost-all  $z \in \mathcal{Z}$  we have:  $k(z|t_0) > 0$ , then also the reverse implication holds, with  $k(x|z, t_0) := \frac{k(x, z|t_0)}{k(z|t_0)}$ .*

$$\forall t \in \mathcal{T} \forall_\mu x, y, z. \quad k(x, y, z|t) = k(x|z, t_0) \cdot k(y, z|t) \quad \longleftarrow \quad X \underset{K(W|T)}{\perp\!\!\!\perp} Y | Z.$$

*Proof.* “ $\implies$ ”: This is clear, as the factorization of the densities provides the needed factorization of the corresponding Markov kernels.

“ $\longleftarrow$ ”: By assumption we have a factorization:

$$K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T),$$



which implies for all  $A \in \mathcal{B}_X$  and  $C \in \mathcal{B}_Z$ :

$$\begin{aligned}
& \int_C \int_A k(x, z|t_0) \mu_X(dx) \mu_Z(dz) \\
&= K(X \in A, Z \in C|T = t_0) \\
&= \int_C Q(X \in A|Z = z) K(Z \in dz|T = t_0) \\
&= \int_C Q(X \in A|Z = z) k(z|t_0) \mu_Z(dz).
\end{aligned}$$

This implies that we have:

$$\forall_{\mu_Z} z \in \mathcal{Z}. \quad \int_A k(x, z|t_0) \mu_X(dx) = Q(X \in A|Z = z) k(z|t_0),$$

which implies, since  $k(z|t_0) > 0$ , that  $k(x|z, t_0) = \frac{k(x, z|t_0)}{k(z|t_0)}$  is a density of  $Q(X|Z)$  up  $\mu_Z$ -null set  $N$ . Since  $K(Z|T) \ll \mu_Z$  this  $N$  is also a  $K(Z|T)$ -null set, and thus  $\mathcal{Y} \times N$  a  $K(Y, Z|T)$ -null set. So for all  $t \in \mathcal{T}$ ,  $A \in \mathcal{B}_X$ ,  $B \in \mathcal{B}_Y$ ,  $C \in \mathcal{B}_Z$  we get:

$$\begin{aligned}
& \int_C \int_B \int_A k(x, y, z|t) \mu_X(dx) \mu_Y(dy) \mu_Z(dz) \\
&= K(X \in A, Y \in B, Z \in C|T = t) \\
&= \int_C \int_B Q(X \in A|Z = z) K(Y \in dy, Z \in dz|T = t) \\
&= \int_C \int_B \int_A k(x|z, t_0) \mu_X(dx) K(Y \in dy, Z \in dz|T = t) \\
&= \int_C \int_B \int_A k(x|z, t_0) \mu_X(dx) k(y, z|t) \mu_Y(dy) \mu_Z(dz) \\
&= \int_C \int_B \int_A k(x|z, t_0) \cdot k(y, z|t) \mu_X(dx) \mu_Y(dy) \mu_Z(dz).
\end{aligned}$$

So the corresponding Markov kernels on the lhs and rhs are the same. This implies that the set:

$$M := \{(x, y, z, t) \mid k(x, y, z|t) \neq k(x|z, t_0) \cdot k(y, z|t)\}$$

is a  $\mu$ -null set in  $\mathcal{B}_X \otimes \mathcal{B}_Y \otimes \mathcal{B}_Z \otimes \mathcal{B}_T$ . So for all  $t \in \mathcal{T}$  and  $\mu$ -almost-all  $x, y, z$  the following equation holds:

$$k(x, y, z|t) = k(x|z, t_0) \cdot k(y, z|t),$$

which implies the claim. □

**Remark 2.5.23** (Conditional independence for random variables). *By Definition/Lemma 2.5.10 we recover the notion of conditional independence for random variables  $X, Y, Z$  with standard measurable spaces as codomains by taking  $\mathcal{T} := \{*\}$ ,  $P(W) := K(W|*)$ :*

$$X \underset{P(W)}{\perp\!\!\!\perp} Y \mid Z \quad \iff \quad \exists Q(X|Z) : P(X, Y, Z) = Q(X|Z) \otimes P(Y, Z).$$

Such a  $Q(X|Z)$  is then a conditional probability distribution of  $P(X, Z)$  conditioned on  $Z$ . In suggestive notations:

$$Q(X|Z) =: P(X|\mathcal{Y}, Z) = P(X|Z).$$

**Lemma 2.5.24** (Conditional independence for deterministic mappings). *Let  $F : \mathcal{T} \rightarrow \mathcal{F}$  and  $H : \mathcal{T} \rightarrow \mathcal{H}$  be measurable mappings, with  $\mathcal{F}$  standard. We now consider them as (deterministic) conditional random variables on the transition probability space  $(\mathcal{W} \times \mathcal{T}, K(W|T))$  via:*

$$\begin{aligned} F : \mathcal{W} \times \mathcal{T} &\rightarrow \mathcal{F}, \\ (w, t) &\mapsto F(t), \\ H : \mathcal{W} \times \mathcal{T} &\rightarrow \mathcal{H}, \\ (w, t) &\mapsto H(t), \end{aligned}$$

which do not depend on the 'probabilistic part'  $\mathcal{W}$  of  $K(W|T)$ . Let  $G : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{G}$  be another conditional random variable.

We write  $F \lesssim H$  if there exists a measurable map  $\varphi : \mathcal{H} \rightarrow \mathcal{F}$  such that:

$$F = \varphi \circ H.$$

Then we have the equivalence:<sup>9</sup>

$$F \lesssim H \iff F \underset{K(W|T)}{\perp\!\!\!\perp} G | H.$$

So  $F$  is a deterministic measurable map of  $H$  iff  $F$  is independent of  $G$  given  $H$ . Note that the first part of the statement is independent of  $G$ .

*Proof.* “ $\implies$ ”: This direction is rather easy. See the later separoid axioms.

“ $\impliedby$ ”: Since  $F$  and  $H$  are deterministic and only dependent on  $T$  we get that:

$$K(F, G, H|T) = \delta(F|T) \otimes \delta(H|T) \otimes K(G|T).$$

By the conditional independence we now have a Markov kernel  $Q(F|H)$  such that we have the factorization:

$$K(F, G, H|T) = Q(F|H) \otimes K(G, H|T) = Q(F|H) \otimes \delta(H|T) \otimes K(G|T).$$

Marginalizing out  $G, H$  and taking  $T = t$  we get from these equations:

$$\delta_{F(t)} = \delta(F|T = t) = Q(F|H(t)),$$

which is a Dirac measure centered at  $F(t)$ . We can now define the mapping:

$$\varphi : H(\mathcal{T}) \rightarrow \mathcal{F}, \quad H(t) \mapsto F(t),$$

---

<sup>9</sup>The full equivalence needs Kuratowski's extension theorem for standard measurable spaces (see [Kec95] 12.2): Any measurable map from a (not necessarily measurable) subset of a measurable space to a standard measurable space extends to a measurable map on the whole space. Alternatively, one could define  $F \lesssim H$  via existence of measurable  $\varphi : H(\mathcal{T}) \rightarrow \mathcal{F}$  such that  $F = \varphi \circ H$ , but this moves problems elsewhere.

which is well-defined, because  $h := H(t_1) = H(t_2)$  implies that  $Q(F|H = h)$  is a Dirac measure centered at  $F(t_1)$  and  $F(t_2)$ , so  $F(t_1) = F(t_2)$ .  $\varphi$  is measurable. Indeed, its composition with  $\delta : \mathcal{F} \rightarrow \mathcal{P}(\mathcal{F})$ ,  $z \mapsto \delta_z$  equals  $Q(F|H)$ , which is measurable. Since  $\mathcal{B}_{\mathcal{F}} = \delta^* \mathcal{B}_{\mathcal{P}(\mathcal{F})}$ , see lemma 2.7.1 2., also  $\varphi$  is measurable. Since  $\mathcal{F}$  is a standard measurable space,  $\varphi$  extends to a measurable mapping  $\varphi : \mathcal{H} \rightarrow \mathcal{F}$  by Kuratowski's extension theorem for standard measurable spaces (see [Kec95] 12.2). Finally, note that we have  $F(t) = \varphi(H(t))$  for all  $(w, t) \in \mathcal{W} \times \mathcal{T}$ , which shows the claim.  $\square$

**Example 2.5.25** (Conditional independence for deterministic mappings). *If for example,  $\mathcal{T} = \mathcal{T}_1 \times \mathcal{T}_2$  and  $T_i : \mathcal{W} \times \mathcal{T}_1 \times \mathcal{T}_2 \rightarrow \mathcal{T}_i$  the canonical projection onto  $\mathcal{T}_i$ , then  $F$  is a function in two variables  $(t_1, t_2)$ . We then have:*

$$F \underset{K(W|T)}{\perp\!\!\!\perp} T_1 | T_2,$$

*if and only if  $F$ —as a function—is only dependent on the argument  $t_2$  (and not on  $t_1$ ).*

Another example of what conditional independence of conditional random variables can encode is the following.

**Remark 2.5.26** (Existence of conditional Markov kernels expressed as conditional independence). *Let  $X, Y$  be conditional random variables on transition probability space  $(\mathcal{W} \times \mathcal{T}, K(W|T))$ . Then we can express the existence of a conditional Markov kernel  $K(X|Y, T)$  as the conditional independence:*

$$X \underset{K(W|T)}{\perp\!\!\!\perp} * | Y, T,$$

*where  $*$  is the constant conditional random variable. Alternatively and equivalently, we could also write:*

$$X \underset{K(W|T)}{\perp\!\!\!\perp} T | Y, T.$$

*Note that for standard measurable spaces  $\mathcal{X}$  and  $\mathcal{Y}$  the above statement always holds. In suggestive symbols:*

$$K(X|\mathcal{X}, Y, T) = K(X|Y, T).$$

**Example 2.5.27** (Certain statistics expressed as conditional independence). *Let  $P(W|\Theta)$  be a statistical model, considered as a Markov kernel  $\mathcal{F} \dashrightarrow \mathcal{W}$ . Let  $X$  and  $Y$  be two conditional random variables w.r.t.  $P(W|\Theta)$ . A statistic of  $X$  is a measurable map  $S : \mathcal{X} \rightarrow \mathcal{S}$ , which we consider as the conditional random variable  $S \lesssim X$  given via:*

$$S : \mathcal{W} \times \mathcal{F} \rightarrow \mathcal{S}, \quad (w, \theta) \mapsto S(X(w, \theta)).$$

1. Ancillarity.  *$S$  is an ancillary statistic of  $X$  w.r.t.  $\Theta$  if and only if:*

$$S \underset{P(W|\Theta)}{\perp\!\!\!\perp} \Theta.$$

*This means that every parameter  $\Theta = \theta$  induces the same distribution for  $S$ :*

$$P(S|\Theta = \theta) = P(S|\emptyset).$$

2. Sufficiency.  $S$  is a sufficient statistic of  $X$  w.r.t.  $\Theta$  if and only if:

$$X \perp\!\!\!\perp_{P(W|\Theta)} \Theta | S.$$

This means that there is a Markov kernel  $P(X|S, \Theta)$  such that:

$$P(X, S|\Theta) = P(X|S, \Theta) \otimes P(S|\Theta).$$

So  $X$  only “interacts” with the parameters  $\Theta$  through  $S$ .

3. Adequacy.  $S$  is an adequate statistic of  $X$  for  $Y$  w.r.t.  $\Theta$  if and only if:

$$X \perp\!\!\!\perp_{P(W|\Theta)} \Theta, Y | S.$$

This means we have a factorization:

$$P(X, Y, S|\Theta) = P(X|\Theta, \mathcal{Y}, S) \otimes P(Y, S|\Theta),$$

for some Markov kernel  $P(X|\Theta, \mathcal{Y}, S)$ . This means that all information of  $X$  about the parameters and labels  $Y$  is fully captured already by  $S$ .

**Theorem 2.5.28.** Let  $(\mathcal{W} \times \mathcal{T}, K(W|T))$  be a transition probability space with Markov kernel:

$$K(W|T) : \mathcal{T} \dashrightarrow \mathcal{W}.$$

Consider conditional random variables  $X, Y, Z$  with common domain  $\mathcal{W} \times \mathcal{T}$  and codomains  $\mathcal{X}, \mathcal{Y}$  and  $\mathcal{Z}$ , resp., and  $T : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{T}$  the canonical projection map. We will write  $P(X|Z) = K(X|\mathcal{X}, Z)$  for a fixed version of the Markov kernel appearing in the conditional independence  $X \perp\!\!\!\perp_{K(W|T)} T | Z$  (only in case it holds).

With these notations, the following are equivalent:

1.  $X \perp\!\!\!\perp_{K(W|T)} Y | Z,$
2.  $X \perp\!\!\!\perp_{K(W|T)} T, Y | Z,$
3.  $X \perp\!\!\!\perp_{K(W|T)} T | Z$  and  $K(X, Y, Z|T) = P(X|Z) \otimes K(Y, Z|T).$
4.  $X \perp\!\!\!\perp_{K(W|T)} T | Z$  and for every  $t \in \mathcal{T}$  we have:  $X_t \perp\!\!\!\perp_{K(W|T=t)} Y_t | Z_t.$

Furthermore, any of those points implies:

$$K(X|\mathcal{X}, Z) = K(X|\mathcal{X}, Z) \quad K(Y, Z|T)\text{-a.s.}$$

and the following:

5. For every probability distribution  $Q(T) \in \mathcal{P}(\mathcal{T})$  we have the conditional independence<sup>10</sup>:

$$X \perp\!\!\!\perp_{K(W|T) \otimes Q(T)} T, Y | Z.$$

*Proof.* 3.  $\implies$  1. is clear by definition.

1.  $\iff$  2.: by 2.5.19.

2.  $\implies$  4.,5.: By assumption we have the factorization:

$$K(X, Y, Z, T|T) = K(X|Z) \otimes K(Y, Z, T|T),$$

for some Markov kernel  $K(X|Z)$ . Via marginalization and multiplication this implies the two equations:

$$\begin{aligned} K(X, Z, T|T) &= K(X|Z) \otimes K(Z, T|T), \\ \underbrace{K(X, Y, Z|T) \otimes Q(T)}_{=:Q(X,Y,Z,T)} &= K(X|Z) \otimes \underbrace{K(Y, Z|T) \otimes Q(T)}_{=:Q(Y,Z,T)}, \end{aligned}$$

for every  $Q(T) \in \mathcal{P}(\mathcal{T})$ . The last equation shows 5.

If we take  $Q(T) = \delta_t$  we get:

$$K(X_t, Y_t, Z_t|T = t) = K(X|Z_t) \otimes K(Y_t, Z_t|T = t).$$

Together with the first of the above equations this shows 4.

4.  $\implies$  3.: By  $X \perp\!\!\!\perp_{K(W|T)} T | Z$  we have:

$$K(X, Z|T) = P(X|Z) \otimes K(Z|T).$$

By the assumption  $X_t \perp\!\!\!\perp_{K(W|T=t)} Y_t | Z_t$ , on the other hand, we have—for each  $t \in \mathcal{T}$  individually—a factorization:

$$K(X, Y, Z|T = t) = Q_t(X|Z) \otimes K(Y, Z|T = t),$$

with a Markov kernel  $Q_t$ , which might depend on  $t \in \mathcal{T}$ , where we suppress the indices  $t$  on all the variables for readability everywhere. Marginalizing out  $Y$  and comparing to the above we then get the two equalities:

$$P(X|Z) \otimes K(Z|T = t) = K(X, Z|T = t) = Q_t(X|Z) \otimes K(Z|T = t).$$

By the essential uniqueness of such a factorization we see that  $P(X \in A|Z)$  only differs from  $Q_t(X \in A|Z)$  on a  $K(Z|T = t)$ -null set. Considered as functions of  $(y, z)$  (by ignoring  $y$ ) they are equal up to  $K(Y, Z|T = t)$ -null set. This means that we can replace  $Q_t(X \in A|Z)$  with  $P(X \in A|Z)$  for every  $A \in \mathcal{B}_X$  and  $t \in \mathcal{T}$ . So we get the equation:

$$K(X, Y, Z|T = t) = P(X|Z) \otimes K(Y, Z|T = t),$$

---

<sup>10</sup>Note that this again implies the second part of point 4:  $X_t \perp\!\!\!\perp Y_t | Z_t$  for every  $t \in \mathcal{T}$ . So the first part of point 4:  $X \perp\!\!\!\perp T | Z$ , can then be seen as the additional obstruction to obtain the “full” conditional independence:  $X \perp\!\!\!\perp Y | Z$ .

for all  $t \in \mathcal{T}$  and thus:

$$K(X, Y, Z|T) = P(X|Z) \otimes K(Y, Z|T).$$

This shows 3. □

**Remark 2.5.29** (Discrete  $\mathcal{T}$ ). *In the setting of theorem 2.5.28, let  $\mathcal{X}$ ,  $\mathcal{Z}$  be standard measurable spaces and  $\mathcal{T}$  be a countable discrete measurable space with any fixed probability distribution  $Q(T)$  that has a strictly positive mass function. Then we get the equivalence:*

$$X \underset{K(W|T) \otimes Q(T)}{\perp\!\!\!\perp} T, Y | Z \iff X \underset{K(W|T)}{\perp\!\!\!\perp} Y | Z.$$

*Proof.* The rhs implies the lhs side by theorem 2.5.28. So, now assume the lhs and put:

$$Q(X, Y, Z, T) := K(X, Y, Z|T) \otimes Q(T).$$

Its marginal is then denoted by  $Q(X, Z)$ . Since  $\mathcal{X}$ ,  $\mathcal{Z}$  are standard measurable spaces we get a (regular) conditional probability distribution  $Q(X|Z)$ , such that  $Q(X, Z) = Q(X|Z) \otimes Q(Z)$ . By the assumed conditional independence we thus have:

$$\begin{aligned} K(X, Y, Z|T) \otimes Q(T) &= Q(X, Y, Z, T) \\ &= Q(X|Z) \otimes Q(Y, Z, T) \\ &= Q(X|Z) \otimes K(Y, Z|T) \otimes Q(T). \end{aligned}$$

Since  $Q(T)$  is strictly positive and conditional Markov kernels are essentially unique we get the sure equality:

$$K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T),$$

which implies the claim. □

**Corollary 2.5.30.** *If  $\mathcal{X}$ ,  $\mathcal{Z}$  are standard measurable spaces then we have the equivalence:*

$$X \underset{K(W|T)}{\perp\!\!\!\perp} Y | Z, T \iff \forall t \in \mathcal{T} : X_t \underset{K(W|T=t)}{\perp\!\!\!\perp} Y_t | Z_t.$$

*Proof.* This directly follows from theorem 2.5.28 4. with  $(Z, T)$  in the role of  $Z$  and remark 2.5.26 to get the first part of 4. In suggestive symbols:

$$K(X|\cancel{T}, Y, Z, T) = K(X|Z, T) \quad K(Z|T)\text{-a.s.}$$

□

### 2.5.4. Example: Linear Gaussian Markov Kernels

**Theorem 2.5.31** (Conditional independence for linear Gaussian conditional random variables). *Let  $\mathcal{T} = \mathbb{R}$ ,  $\mathcal{X} = \mathbb{R}^{d_x}$ ,  $\mathcal{Y} = \mathbb{R}^{d_y}$  and  $\mathcal{Z} = \mathbb{R}^{d_z}$ . Consider a linear Gaussian Markov kernel  $P(X, Y, Z|T)$ , which is given by a density of the form:*

$$p(x, y, z|t) = \mathcal{N} \left( \begin{bmatrix} x \\ y \\ z \end{bmatrix} \middle| \begin{bmatrix} \Gamma_X \\ \Gamma_Y \\ \Gamma_Z \end{bmatrix} \cdot t + \begin{bmatrix} \gamma_X \\ \gamma_Y \\ \gamma_Z \end{bmatrix}, \begin{bmatrix} \Sigma_{X,X} & \Sigma_{X,Y} & \Sigma_{X,Z} \\ \Sigma_{Y,X} & \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,X} & \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix} \right).$$

Then we have the following equivalence:

$$X \perp\!\!\!\perp_{P(X,Y,Z|T)} Y | Z \quad \iff \quad \Sigma_{X,Y} = \Sigma_{X,Z} \Sigma_{Z,Z}^{-1} \Sigma_{Z,Y} \quad \wedge \quad \Gamma_X = \Sigma_{X,Z} \Sigma_{Z,Z}^{-1} \Gamma_Z.$$

If this is the case then the Markov kernel  $Q(X|Z)$  coming from the conditional independence:

$$P(X, Y, Z|T) = Q(X|Z) \otimes P(Y, Z|T),$$

is also a linear Gaussian Markov kernel with density:

$$\begin{aligned} q(x|z) &= \mathcal{N}(x | \mu_{X|Z}(z), \Sigma_{X,X|Z}), \\ \mu_{X|Z}(z) &:= \gamma_X + \Sigma_{X,Z} \Sigma_{Z,Z}^{-1} (z - \gamma_Z), \\ \Sigma_{X,X|Z} &:= \Sigma_{X,X} - \Sigma_{X,Z} \Sigma_{Z,Z}^{-1} \Sigma_{Z,X}, \end{aligned}$$

which coincides with the usual marginal conditional for  $t = 0$ , i.e.:

$$Q(X|Z = z) = P(X|Z = z, T = 0).$$

So, we also get the equivalence:

$$X \perp\!\!\!\perp_{P(X,Y,Z|T)} Y | Z \quad \iff \quad P(X, Y, Z|T) = P(X|Z, T = 0) \otimes P(Y, Z|T).$$

*Proof.* First note that in general the conditional  $P(X|Y, Z, T)$  is also a linear Gaussian Markov kernel and of the form:

$$p(x|y, z, t) = \mathcal{N}(x | \mu_{X|Y,Z,T}(y, z, t), \Sigma_{X|Y,Z,T}),$$

with the following abbreviation for the covariance matrix:

$$\Sigma_{X|Y,Z,T} := \Sigma_{X,X} - \begin{bmatrix} \Sigma_{X,Y} & \Sigma_{X,Z} \end{bmatrix} \begin{bmatrix} \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix}^{-1} \begin{bmatrix} \Sigma_{Y,X} \\ \Sigma_{Z,X} \end{bmatrix},$$

and the following abbreviation for the mean:

$$\begin{aligned} &\mu_{X|Y,Z,T}(y, z, t) \\ &:= (\Gamma_X \cdot t + \gamma_X) + \begin{bmatrix} \Sigma_{X,Y} & \Sigma_{X,Z} \end{bmatrix} \begin{bmatrix} \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix}^{-1} \left( \begin{bmatrix} y \\ z \end{bmatrix} - \begin{bmatrix} \Gamma_Y \cdot t + \gamma_Y \\ \Gamma_Z \cdot t + \gamma_Z \end{bmatrix} \right) \\ &= \begin{bmatrix} \Sigma_{X,Y} & \Sigma_{X,Z} \end{bmatrix} \begin{bmatrix} \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix}^{-1} \begin{bmatrix} y \\ z \end{bmatrix} \\ &\quad + (\Gamma_X \cdot t + \gamma_X) - \begin{bmatrix} \Sigma_{X,Y} & \Sigma_{X,Z} \end{bmatrix} \begin{bmatrix} \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix}^{-1} \begin{bmatrix} \Gamma_Y \cdot t + \gamma_Y \\ \Gamma_Z \cdot t + \gamma_Z \end{bmatrix}. \end{aligned}$$

We now want to investigate under which conditions we get the conditional independence:

$$X \perp\!\!\!\perp_{P(X,Y,Z|T)} Y \mid Z.$$

Note that in this case the conditional independence is equivalent to the statement that the conditional density  $p(x|y, z, t)$  is not dependent on the arguments  $y$  and  $t$ .

Let us first investigate the first term involving  $y$ :

$$[\Sigma_{X,Y} \quad \Sigma_{X,Z}] \begin{bmatrix} \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix}^{-1} \begin{bmatrix} y \\ z \end{bmatrix}.$$

Note that we can use the following formula for the  $(2 \times 2)$ -block inverse:

$$\begin{bmatrix} \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix}^{-1} = \begin{bmatrix} (\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1} & -(\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1}\Sigma_{Y,Z}\Sigma_{Z,Z}^{-1} \\ -\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y}(\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1} & \Sigma_{Z,Z}^{-1} + \Sigma_{Z,Z}^{-1}\Sigma_{Z,Y}(\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1}\Sigma_{Y,Z}\Sigma_{Z,Z}^{-1} \end{bmatrix}$$

This leads us to require that:

$$[\Sigma_{X,Y} \quad \Sigma_{X,Z}] \begin{bmatrix} (\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1} \\ -\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y}(\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1} \end{bmatrix} = 0,$$

which is equivalent to:

$$(\Sigma_{X,Y} - \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})(\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1} = 0,$$

which can be further simplified, by multiplying with the inverse of the inverse, to:

$$\Sigma_{X,Y} - \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y} = 0.$$

We also need that the mean of the conditional is not dependent on  $t$ , which leads to the following condition coming from the second term:

$$\begin{aligned} 0 &= \Gamma_X - [\Sigma_{X,Y} \quad \Sigma_{X,Z}] \begin{bmatrix} \Sigma_{Y,Y} & \Sigma_{Y,Z} \\ \Sigma_{Z,Y} & \Sigma_{Z,Z} \end{bmatrix}^{-1} \begin{bmatrix} \Gamma_Y \\ \Gamma_Z \end{bmatrix} \\ &= \Gamma_X - (\Sigma_{X,Y} - \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})(\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1}\Gamma_Y \\ &\quad - \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Gamma_Z + (\Sigma_{X,Y} - \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})(\Sigma_{Y,Y} - \Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y})^{-1}\Sigma_{Y,Z}\Sigma_{Z,Z}^{-1}\Gamma_Z \\ &= \Gamma_X - \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Gamma_Z, \end{aligned}$$

where we made repeated use of the condition:  $\Sigma_{X,Y} - \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y} = 0$ .

This leads us to the following equivalence for linear Gaussian Markov kernels:

$$X \perp\!\!\!\perp_{P(X,Y,Z|T)} Y \mid Z \iff \Sigma_{X,Y} = \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Sigma_{Z,Y} \quad \wedge \quad \Gamma_X = \Sigma_{X,Z}\Sigma_{Z,Z}^{-1}\Gamma_Z.$$



If this is the case then the Markov kernel  $Q(X|Z)$  coming from the conditional independence:

$$P(X, Y, Z|T) = Q(X|Z) \otimes P(Y, Z|T),$$

is also a linear Gaussian Markov kernel with density:

$$\begin{aligned} q(x|z) &= \mathcal{N}(x \mid \mu_{X|Z}(z), \Sigma_{X,X|Z}), \\ \mu_{X|Z}(z) &:= \gamma_X + \Sigma_{X,Z} \Sigma_{Z,Z}^{-1} (z - \gamma_Z), \\ \Sigma_{X,X|Z} &:= \Sigma_{X,X} - \Sigma_{X,Z} \Sigma_{Z,Z}^{-1} \Sigma_{Z,X}, \end{aligned}$$

which is the usual marginal conditional for  $t = 0$ . □

## 2.6. Separoid Axioms for Conditional Independence

The following asymmetric separoid axioms for conditional independence are a generalization of the symmetric separoid axioms due to A.P. Dawid [Daw01] and the similar graphoid axioms due to J. Pearl and A. Paz [PP85].

**Definition/Theorem 2.6.1** ((Asymmetric) separoid axioms for conditional independence). *Let  $(\mathcal{W} \times \mathcal{T}, K(W|T))$  be a transition probability space with Markov kernel:*

$$K(W|T) : \mathcal{T} \dashrightarrow \mathcal{W}.$$

*Consider conditional random variables  $X, Y, Z, U$  with common domain  $\mathcal{W} \times \mathcal{T}$  and standard measurable spaces  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}, \mathcal{U}$ , resp., as codomains. Let  $T : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{T}$  be the canonical projection and  $*$  the constant conditional random variable.*

*We write  $U \lesssim X$  if there exists a measurable function  $G : \mathcal{X} \rightarrow \mathcal{U}$  such that  $U = G \circ X$ . Then the ternary relation  $\perp\!\!\!\perp = \perp\!\!\!\perp_{K(W|T)}$  satisfies the following rules:*

a) *Extended Left Redundancy:*

$$U \lesssim X \implies U \perp\!\!\!\perp Y \mid X.$$

b) *T-Restricted Right Redundancy:<sup>11</sup>*

$$X \perp\!\!\!\perp * \mid Z, T \text{ always holds.}$$

c) *T-Inverted Right Decomposition:*

$$X \perp\!\!\!\perp Y \mid Z \implies X \perp\!\!\!\perp T, Y \mid Z.$$

d) *Left Decomposition:*

$$X, U \perp\!\!\!\perp Y \mid Z \implies U \perp\!\!\!\perp Y \mid Z.$$

e) *Right Decomposition:*

---

<sup>11</sup> T-Restricted Right Redundancy, Left Weak Union and Symmetry need the existence of conditional Markov kernels. That is the reason we assumed standard measurable spaces.

$$X \perp\!\!\!\perp Y, U \mid Z \implies X \perp\!\!\!\perp U \mid Z.$$

f) *Left Weak Union:*<sup>11</sup>

$$X, U \perp\!\!\!\perp Y \mid Z \implies X \perp\!\!\!\perp Y \mid U, Z.$$

g) *Right Weak Union:*

$$X \perp\!\!\!\perp Y, U \mid Z \implies X \perp\!\!\!\perp Y \mid U, Z.$$

h) *Left Contraction:*

$$(X \perp\!\!\!\perp Y \mid U, Z) \wedge (U \perp\!\!\!\perp Y \mid Z) \implies X, U \perp\!\!\!\perp Y \mid Z.$$

i) *Right Contraction:*

$$(X \perp\!\!\!\perp Y \mid U, Z) \wedge (X \perp\!\!\!\perp U \mid Z) \implies X \perp\!\!\!\perp Y, U \mid Z.$$

j) *Right Cross Contraction:*

$$(X \perp\!\!\!\perp Y \mid U, Z) \wedge (U \perp\!\!\!\perp X \mid Z) \implies X \perp\!\!\!\perp Y, U \mid Z.$$

k) *Flipped Left Cross Contraction:*

$$(X \perp\!\!\!\perp Y \mid U, Z) \wedge (Y \perp\!\!\!\perp U \mid Z) \implies Y \perp\!\!\!\perp X, U \mid Z.$$

*In particular, we have the equivalences:*

$$(X \perp\!\!\!\perp Y, U \mid Z) \iff (X \perp\!\!\!\perp Y \mid U, Z) \wedge (X \perp\!\!\!\perp U \mid Z),$$

$$(X, U \perp\!\!\!\perp Y \mid Z) \iff (X \perp\!\!\!\perp Y \mid U, Z) \wedge (U \perp\!\!\!\perp Y \mid Z).$$

*We also get:*

l) *T-Restricted Symmetry:*<sup>11</sup>

$$X \perp\!\!\!\perp Y \mid Z, T \implies Y \perp\!\!\!\perp X \mid Z, T.$$

*In the special case of  $\mathcal{T} = * = \{*\}$ , the one-point space, (i.e. in the case of probability distributions and random variables mapping to standard measurable spaces) we thus have (unrestricted) Symmetry.*

## Proofs - Separoid Axioms for Conditional Independence

In the following let  $(\mathcal{W} \times \mathcal{T}, K(W|T))$  be a transition probability space with Markov kernel:

$$K(W|T) : \mathcal{T} \dashrightarrow \mathcal{W},$$

and conditional random variables  $X, Y, Z, U$  with common domain  $\mathcal{W} \times \mathcal{T}$  and measurable spaces  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}, \mathcal{U}$ , resp., as codomains. We indicate when we need to assume standard measurable spaces.

We will use  $T : \mathcal{W} \times \mathcal{T} \rightarrow \mathcal{T}$  to denote the canonical projection and  $*$  to denote the constant conditional random variable.

Recall that we write  $U \lesssim X$  if there exists a measurable function  $\varphi : \mathcal{X} \rightarrow \mathcal{U}$  such that  $U = \varphi \circ X$ .

Recall that for proving:

$$X \perp\!\!\!\perp_{K(W|T)} Y | Z,$$

we need to find/construct a Markov kernel  $Q(X|Z)$  such that:

$$K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T),$$

which is equivalent to:

For all  $t \in \mathcal{T}$  and all measurable  $A \subseteq \mathcal{X}, B \subseteq \mathcal{Y}, C \subseteq \mathcal{Z}$  we have the equation:

$$K(X \in A, Y \in B, Z \in C | T = t) = \int_C \int_B Q(X \in A | Z = z) K(Y \in dy, Z \in dz | T = t).$$

We abbreviate  $\perp\!\!\!\perp := \perp\!\!\!\perp_{K(W|T)}$  in the following.

**Lemma 2.6.2** (Extended Left Redundancy).

$$U \lesssim X \implies U \perp\!\!\!\perp Y | X.$$

*Proof.* If  $U = \varphi(X)$  put  $Q(U \in D | X = x) := \delta_\varphi(U \in D | X = x) := \mathbb{1}_D(\varphi(x))$  for  $D \subseteq \mathcal{U}$ . Then we get:

$$\begin{aligned} & \int_C \int_B Q(U \in D | X = x) K(Y \in dy, X \in dx | T = t) \\ &= \int_C \int_B \delta_\varphi(U \in D | X = x) K(Y \in dy, X \in dx | T = t) \\ &= \int_C \int_B \mathbb{1}_{\varphi^{-1}(D)}(x) K(Y \in dy, X \in dx | T = t) \\ &= K(Y \in B, X \in C \cap \varphi^{-1}(D) | T = t) \\ &= K(Y \in B, X \in C, \varphi(X) \in D | T = t) \\ &= K(U \in D, Y \in B, X \in C | T = t). \end{aligned}$$

This shows the claim. In suggestive symbols:

$$K(U | \cancel{\mathcal{T}}, \mathcal{Y}, X) = \delta_\varphi(U | X).$$

□

**Lemma 2.6.3** (*T-Restricted Right Redundancy*). *Let  $\mathcal{X}$  and  $\mathcal{Z}$  be standard measurable spaces. Then:*

$$X \perp\!\!\!\perp * | Z, T \quad \text{holds.}$$

*Proof.* Because  $\mathcal{X}$  and  $\mathcal{Z}$  are standard measurable spaces we have a factorization:

$$K(X, *, Z, T|T) = K(X|Z, T) \otimes K(*, Z, T|T).$$

with the conditional Markov kernel  $K(X|Z, T)$  of  $K(X, Z|T)$  (via theorem 2.4.16). This already shows the claim. In suggestive symbols:

$$K(X|_{\cancel{*}}, \cancel{\mathcal{T}}, Z, T) = K(X|Z, T).$$

□

**Lemma 2.6.4** (*T-Inverted Right Decomposition*).

$$X \perp\!\!\!\perp Y | Z \implies X \perp\!\!\!\perp T, Y | Z.$$

*Proof.* By assumption we have:

$$K(X, Y, Z|T) = Q(X|Z) \otimes K(Y, Z|T).$$

Multiplying both sides with  $\delta(T|T)$  we get:

$$K(X, Y, Z, T|T) = Q(X|Z) \otimes K(T, Y, Z|T).$$

This shows the claim using the same  $Q(X|Z)$ . In suggestive symbols:

$$K(X|_{\cancel{T}}, \cancel{\mathcal{T}}, Y, Z) = K(X|Z, Y, Z).$$

□

**Lemma 2.6.5** (*Left Decomposition*).

$$X, U \perp\!\!\!\perp Y | Z \implies U \perp\!\!\!\perp Y | Z.$$

*Proof.* Let  $Q(X, U|Z)$  be given from the left conditional independence. Then we have:

$$K(X, U, Y, Z|T) = Q(X, U|Z) \otimes K(Y, Z|T).$$

Marginalizing out  $X$  gives:

$$K(U, Y, Z|T) = Q(U|Z) \otimes K(Y, Z|T).$$

This shows the claim. In suggestive symbols:

$$K(U|_{\cancel{X}}, \cancel{\mathcal{X}}, Y, Z) = K(U|Z, Y, Z).$$

□

**Lemma 2.6.6** (Right Decomposition).

$$X \perp\!\!\!\perp Y, U \mid Z \implies X \perp\!\!\!\perp U \mid Z.$$

*Proof.* Let  $Q(X|Z)$  be given from the left conditional independence. We then have:

$$K(X, U, Y, Z|T) = Q(X|Z) \otimes K(Y, U, Z|T).$$

Marginalizing out  $Y$  gives:

$$K(X, U, Z|T) = Q(X|Z) \otimes K(U, Z|T).$$

This shows the claim. In suggestive symbols:

$$K(X|\cancel{T}, \mathcal{U}, Z) = K(X|\cancel{T}, Y, \mathcal{U}, Z).$$

□

**Lemma 2.6.7** (Left Weak Union). *Let  $\mathcal{X}$  and  $\mathcal{U}$  be standard measurable spaces. Then:*

$$X, U \perp\!\!\!\perp Y \mid Z \implies X \perp\!\!\!\perp Y \mid U, Z.$$

*Proof.* By assumption we have:

$$K(X, U, Y, Z|T) = Q(X, U|Z) \otimes K(Y, Z|T),$$

for some Markov kernel  $Q(X, U|Z)$ . If we marginalize out  $X$  we get:

$$K(U, Y, Z|T) = Q(U|Z) \otimes K(Y, Z|T).$$

Because  $\mathcal{X}$  and  $\mathcal{U}$  are standard measurable spaces we have a factorization:

$$Q(X, U|Z) = Q(X|U, Z) \otimes Q(U|Z).$$

with the conditional Markov kernel  $Q(X|U, Z)$  (via theorem 2.4.16).

Putting these equations together we get:

$$\begin{aligned} K(X, U, Y, Z|T) &= Q(X, U|Z) \otimes K(Y, Z|T) \\ &= Q(X|U, Z) \otimes Q(U|Z) \otimes K(Y, Z|T) \\ &= Q(X|U, Z) \otimes K(U, Y, Z|T). \end{aligned}$$

In suggestive symbols, this means that:  $K(X|\cancel{T}, \mathcal{Y}, U, Z)$  is the conditional of  $K(X, U|\cancel{T}, \mathcal{Y}, Z)$ .

□

**Lemma 2.6.8** (Right Weak Union).

$$X \perp\!\!\!\perp Y, U \mid Z \implies X \perp\!\!\!\perp Y \mid U, Z.$$

*Proof.* We have the factorization:

$$K(X, Y, U, Z|T) = Q(X|Z) \otimes K(Y, U, Z|T),$$

with some Markov kernel  $Q(X|Z)$ . If we view  $Q(X|Z)$  as a function in  $(u, z)$  via:

$$(u, z) \mapsto Q(X|Z = z),$$

by just ignoring the argument  $u$  then the claim follows from the same factorization above.

In suggestive symbols:

$$K(X|\cancel{T}, Y, U, Z) = K(X|\cancel{T}, Y, \cancel{U}, Z).$$

□

**Lemma 2.6.9** (Left Contraction).

$$(X \perp\!\!\!\perp Y | U, Z) \wedge (U \perp\!\!\!\perp Y | Z) \implies X, U \perp\!\!\!\perp Y | Z.$$

*Proof.* By assumption we have the two factorizations:

$$\begin{aligned} K(X, Y, U, Z|T) &= Q(X|U, Z) \otimes K(Y, U, Z|T), \\ K(Y, U, Z|T) &= P(U|Z) \otimes K(Y, Z|T), \end{aligned}$$

with some Markov kernels  $Q(X|U, Z)$ ,  $P(U|Z)$ . Putting these equations together using  $Q(X|U, Z) \otimes P(U|Z)$  we get:

$$K(X, Y, U, Z|T) = (Q(X|U, Z) \otimes P(U|Z)) \otimes K(Y, Z|T).$$

In suggestive symbols:

$$K(X, U|\cancel{T}, Y, Z) = K(X|\cancel{T}, Y, U, Z) \otimes K(U|\cancel{T}, Y, Z).$$

□

**Lemma 2.6.10** (Right Contraction).

$$(X \perp\!\!\!\perp Y | U, Z) \wedge (X \perp\!\!\!\perp U | Z) \implies X \perp\!\!\!\perp Y, U | Z.$$

*Proof.* By assumption we have the two factorizations:

$$\begin{aligned} K(X, Y, U, Z|T) &= Q(X|U, Z) \otimes K(Y, U, Z|T), \\ K(X, U, Z|T) &= P(X|Z) \otimes K(U, Z|T), \end{aligned}$$

with some Markov kernels  $Q(X|U, Z)$ ,  $P(X|Z)$ .

Marginalizing out  $Y$  we get the equalities:

$$\begin{aligned} K(X, U, Z|T) &= Q(X|U, Z) \otimes K(U, Z|T), \\ K(X, U, Z|T) &= P(X|Z) \otimes K(U, Z|T). \end{aligned}$$

By the essential uniqueness (see lemma 2.4.22) of such factorization we get that for every  $A \in \mathcal{B}_X$ :

$$Q(X \in A|U, Z) = P(X \in A|Z) \quad K(U, Z|T)\text{-a.s.}$$

The same equation then holds also  $K(Y, U, Z|T)$ -a.s. (by ignoring argument  $y$ ). Plugging that back into the first equation gives:

$$K(X, Y, U, Z|T) = P(X|Z) \otimes K(Y, U, Z|T).$$

In suggestive symbols:

$$K(X|\underline{T}, \underline{Y}, \underline{U}, Z) = K(X|\underline{T}, \underline{Y}, U, Z) = K(X|\underline{T}, \underline{U}, Z) \quad \text{a.s.}$$

□

**Lemma 2.6.11** (Right Cross Contraction).

$$(X \perp\!\!\!\perp Y | U, Z) \wedge (U \perp\!\!\!\perp X | Z) \implies X \perp\!\!\!\perp Y, U | Z.$$

*Proof.* By assumption we have the two factorizations:

$$K(X, Y, U, Z|T) = Q(X|U, Z) \otimes K(Y, U, Z|T), \quad (3)$$

$$K(X, U, Z|T) = P(U|Z) \otimes K(X, Z|T), \quad (4)$$

with some Markov kernels  $Q(X|U, Z)$ ,  $P(U|Z)$ .

We then define the Markov kernel:

$$R(X, U|Z) := Q(X|U, Z) \otimes P(U|Z). \quad (5)$$

We will now show that its marginal:

$$R(X|Z) = Q(X|U, Z) \circ P(U|Z). \quad (6)$$

will satisfy the claim.

If we marginalize out  $Y$  from equation 3 we get:

$$K(X, U, Z|T) = Q(X|U, Z) \otimes K(U, Z|T). \quad (7)$$

Equating equations 4 and 7 gives:

$$P(U|Z) \otimes K(X, Z|T) = K(X, U, Z|T) = Q(X|U, Z) \otimes K(U, Z|T). \quad (8)$$

Marginalizing out  $X$  in equation 8 on both sides gives:

$$K(U, Z|T) = P(U|Z) \otimes K(Z|T). \quad (9)$$

If we now plug equation 9 into 7 then we get:

$$K(X, U, Z|T) = Q(X|U, Z) \otimes P(U|Z) \otimes K(Z|T) \quad (10)$$

$$\stackrel{5}{=} R(X, U|Z) \otimes K(Z|T). \quad (11)$$

If we marginalize out  $U$  in equation 11 and use equation 6 we arrive at:

$$K(X, Z|T) = R(X|Z) \otimes K(Z|T). \quad (12)$$

We now get:

$$Q(X|U, Z) \otimes K(U, Z|T) \stackrel{7}{=} K(X, U, Z|T) \quad (13)$$

$$\stackrel{4}{=} P(U|Z) \otimes K(X, Z|T) \quad (14)$$

$$\stackrel{12}{=} P(U|Z) \otimes R(X|Z) \otimes K(Z|T) \quad (15)$$

$$= R(X|Z) \otimes P(U|Z) \otimes K(Z|T) \quad (16)$$

$$\stackrel{9}{=} R(X|Z) \otimes K(U, Z|T). \quad (17)$$

By the essential uniqueness (see lemma 2.4.22) of such a factorization we get that for every  $A \in \mathcal{B}_X$ :

$$Q(X \in A|U, Z) = R(X \in A|Z) \quad K(U, Z|T)\text{-a.s.} \quad (18)$$

The same equation then holds also  $K(Y, U, Z|T)$ -a.s. (by ignoring the non-occurring argument  $y$ ). Plugging equation 18 back into the equation 3 we get:

$$K(X, Y, U, Z|T) = Q(X|U, Z) \otimes K(Y, U, Z|T), \quad (19)$$

$$= R(X|Z) \otimes K(Y, U, Z|T). \quad (20)$$

This shows the claim.

In suggestive symbols:

$$K(X|\cancel{T}, \cancel{Y}, U, Z) = K(X|\cancel{T}, \cancel{Y}, U, Z) \circ K(U|\cancel{T}, \cancel{X}, Z).$$

□

**Lemma 2.6.12** (Flipped Left Cross Contraction).

$$(X \perp\!\!\!\perp Y | U, Z) \wedge (Y \perp\!\!\!\perp U | Z) \implies Y \perp\!\!\!\perp X, U | Z.$$

*Proof.* By assumption we have the two factorizations:

$$K(X, Y, U, Z|T) = Q(X|U, Z) \otimes K(Y, U, Z|T),$$

$$K(Y, U, Z|T) = P(Y|Z) \otimes K(U, Z|T),$$

with some Markov kernels  $Q(X|U, Z)$ ,  $P(Y|Z)$ .

Marginalizing out  $Y$  in the first equation we get the equality:

$$K(X, U, Z|T) = Q(X|U, Z) \otimes K(U, Z|T).$$



Plugging all three equations into each other we get:

$$\begin{aligned}
K(X, Y, U, Z|T) &= Q(X|U, Z) \otimes K(Y, U, Z|T) \\
&= Q(X|U, Z) \otimes P(Y|Z) \otimes K(U, Z|T) \\
&= P(Y|Z) \otimes Q(X|U, Z) \otimes K(U, Z|T) \\
&= P(Y|Z) \otimes K(X, U, Z|T).
\end{aligned}$$

In suggestive symbols:

$$K(Y|\cancel{T}, \cancel{X}, U, Z) = K(Y|\cancel{T}, U, Z).$$

□

**Lemma 2.6.13** (*T*-Restricted Symmetry). *Let  $\mathcal{Y}$  and  $\mathcal{Z}$  be standard measurable spaces. Then:*

$$X \perp\!\!\!\perp Y | Z, T \implies Y \perp\!\!\!\perp X | Z, T.$$

*Proof.* This follows from Flipped Left Cross Contraction with  $U = *$  and  $(Z, T)$  for  $Z$ :

$$(X \perp\!\!\!\perp Y | Z, T) \wedge (Y \perp\!\!\!\perp * | Z, T) \implies Y \perp\!\!\!\perp X | Z, T,$$

together with *T*-Restricted Right Redundancy:

$$Y \perp\!\!\!\perp * | Z, T.$$

In suggestive symbols:

$$K(Y|\cancel{*}, \cancel{T}, X, Z) = K(Y|\cancel{T}, X, Z).$$

□

## 2.7. Markov Kernels from Deterministic Mappings

**Lemma 2.7.1.** *Let  $\mathcal{X}$ ,  $\mathcal{Y}$ ,  $\mathcal{Z}$  be measurable spaces.*

1. *If  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is measurable then the induced map:*

$$f_* : \mathcal{P}(\mathcal{X}) \rightarrow \mathcal{P}(\mathcal{Y}), \quad P \mapsto f_*P = (B \mapsto P(f^{-1}(B))),$$

*is measurable as well.*

2. *The map:*

$$\delta : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{X}), \quad x \mapsto \delta_x = (A \mapsto \mathbb{1}_A(x)),$$

*is measurable and  $\delta^*\mathcal{B}_{\mathcal{P}(\mathcal{X})} = \mathcal{B}_{\mathcal{X}}$ .  $\delta$  is injective iff  $\mathcal{B}_{\mathcal{X}}$  separates points.*

3. *The map:*

$$\begin{aligned}
\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y}) &\rightarrow \mathcal{P}(\mathcal{X} \times \mathcal{Y}), \\
(P, Q) &\mapsto P \otimes Q,
\end{aligned}$$

*is measurable.*

4. If  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Z}$  is measurable then the map:

$$\begin{aligned} \mathcal{P}(\mathcal{X}) \times \mathcal{Y} &\rightarrow \mathcal{P}(\mathcal{Z}), \\ (P, y) &\mapsto g_*(P \otimes \delta_y) \\ &= (C \mapsto P(\{x \in \mathcal{X} \mid g(x, y) \in C\})) \end{aligned}$$

is measurable as well.

**Remark 2.7.2.** Let  $f : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathcal{X}$  be measurable and  $P(Y) \in \mathcal{P}(\mathcal{Y})$  a fixed probability distribution. Then the map:

$$K(X|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}, \quad (A, z) \mapsto P(f(Y, z) \in A) =: K(X \in A|Z = z)$$

is a Markov kernel.

**Theorem 2.7.3.** Let  $\mathcal{Z}$  be any measurable space and  $\mathcal{X} = \bar{\mathbb{R}} = [-\infty, \infty]$ . Let  $K(X|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}$  be a Markov kernel,  $P(E)$  be the uniform distribution on  $\mathcal{E} := [0, 1]$  and:

$$R(e|z) := \inf \{ \tilde{x} \in \mathcal{X} \mid K(X \leq \tilde{x}|z) \geq e \},$$

the (conditional) quantile function (a.k.a. inverse cumulative distribution function) of  $K(X|Z)$ . Then we can write  $K(X|Z)$  as the push-forward:

$$K(X|Z) = \delta(R|E, Z) \circ P(E).$$

More explicitly, for  $A \in \mathcal{B}_{\mathcal{X}}$  and  $z \in \mathcal{Z}$  we have:

$$K(X \in A|Z = z) = P(R(E|z) \in A).$$

*Proof.* We only need to check the last equation for  $A = [-\infty, x]$  and  $x \in \bar{\mathbb{R}}$ . We then use the following equivalence for  $x \in \bar{\mathbb{R}}$ ,  $z \in \mathcal{Z}$  and  $e \in [0, 1]$ , see Lemma 2.7.8:

$$R(e|z) \leq x \quad \iff \quad e \leq F(x|z),$$

where  $F$  is the conditional cumulative distribution function of  $K(X|Z)$ . So we get:

$$P(R(E|z) \leq x) = P(E \leq F(x|z)) = F(x|z) := K(X \leq x|Z = z).$$

The equality in the middle holds because  $E$  is uniformly distributed. This shows the claim.  $\square$

**Remark 2.7.4.** Let  $\mathcal{Z}$  be any measurable space and  $\mathcal{X}$  be a standard measurable space with a fixed embedding  $\iota : \mathcal{X} \hookrightarrow \bar{\mathbb{R}} = [-\infty, \infty]$  onto a Borel subset, which always exists, and  $K(X|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}$  a Markov kernel. Then the push-forward Markov kernel:

$$K(\iota X|Z) : \mathcal{Z} \xrightarrow{K(X|Z)} \mathcal{P}(\mathcal{X}) \xrightarrow{\iota_*} \mathcal{P}(\bar{\mathbb{R}}), \quad (A, z) \mapsto K(X \in \iota^{-1}(A)|Z = z),$$

satisfies the condition of Theorem 2.7.3. So with those notations we get for all  $A \in \mathcal{B}_{\bar{\mathbb{R}}}$  and  $z \in \mathcal{Z}$ :

$$K(\iota X \in A|Z = z) = P(R(E|z) \in A).$$

Since  $\iota(\mathcal{X}) \in \mathcal{B}_{\mathbb{R}}$  we get for all  $z \in \mathcal{Z}$ :

$$0 = K(\iota X \in \iota(\mathcal{X})^c | Z = z) = P(R(E|z) \in \iota(\mathcal{X})^c).$$

Since  $R$  is measurable we get that:

$$D := \{(e, z) \in \mathcal{E} \times \mathcal{Z} \mid R(e|z) \in \iota(\mathcal{X})\} \in \mathcal{B}_{\mathcal{E}} \otimes \mathcal{B}_{\mathcal{Z}},$$

with  $P(E \in D_z^c) = 0$  for all  $z \in \mathcal{Z}$ . We can then measurably adjust  $R$  to get a measurable map:

$$\tilde{R}: \mathcal{E} \times \mathcal{Z} \rightarrow \mathcal{X}, \quad \tilde{R}(e|z) := \iota^{-1}(R(e|z)),$$

for  $(e, z) \in D$  and  $\tilde{R}(e|z) := \tilde{x}$  for  $(e, z) \in D^c$  and a fixed point  $\tilde{x} \in \mathcal{X}$ . With this adjustment we then get for all  $A \in \mathcal{B}_{\mathcal{X}}$  and  $z \in \mathcal{Z}$ :

$$\begin{aligned} K(X \in A | Z = z) &= K(\iota X \in \iota(A) | Z = z) \\ &= P(R(E|z) \in \iota(A)) \\ &= P(\iota(\tilde{R}(E|z)) \in \iota(A)) \\ &= P(\tilde{R}(E|z) \in A), \end{aligned}$$

or in short:

$$K(X|Z) = \delta(\tilde{R}|E, Z) \circ P(E).$$

In other words, Theorem 2.7.3 holds (with those slight adjustments) for all standard measurable spaces  $\mathcal{X}$  as well.

In terms of random variables the theorem above states that every distribution  $Q$  can be generated by the uniform one  $U[0, 1]$  and a deterministic map. The theorem below strengthens this claim. It says that every conditional random variable  $X$  can be represented in terms of a uniformly distributed random variable  $E$  and a measurable map. In short, the above is about 'in distribution' and the one below about 'almost-surely' statements.

**Theorem 2.7.5.** Let  $\mathcal{X} := \bar{\mathbb{R}}$ ,  $\mathcal{U} := [0, 1]$  and  $\mathcal{Z}$  any measurable space. Let  $X$ ,  $U$  and  $Z$  be conditional random variables taking values in  $\mathcal{X}$ ,  $\mathcal{U}$  and  $\mathcal{Z}$ , resp., such that:

$$U \underset{K(U, X|Z)}{\perp\!\!\!\perp} X, Z,$$

with  $K(U|X, Z)$  the uniform distribution on  $[0, 1]$ . Define the interpolated (conditional) cumulative distribution function and its corresponding quantile function via:

$$\begin{aligned} F(x; u|z) &:= K(X < x | Z = z) + u \cdot K(X = x | Z = z), \\ R(e|z) &:= \inf \{\tilde{x} \in \mathcal{X} \mid F(\tilde{x}; 1|z) \geq e\}, \end{aligned}$$

and the conditional random variable  $E := F(X; U|Z)$ . Then we have the (conditional) independence:

$$E \underset{K(U, X|Z)}{\perp\!\!\!\perp} Z,$$

with  $K(E|Z)$  the uniform distribution on  $[0, 1]$ , and:

$$X = R(E|Z) \quad K(U, X|Z)\text{-a.s.}$$

*Proof.* After the (joint) measurabilities of  $F$  and  $R$  are checked the statement directly follows from Lemma 2.7.10 by applying it for every  $z$  separately.  $\square$

**Remark 2.7.6.** *With a similar argument as used in Remark 2.7.4 we can in Theorem 2.7.5 replace  $\mathcal{X}$  by any standard measurable space. We then use  $E := F(\iota X; U|Z)$  to get the conditional independence:*

$$E \perp\!\!\!\perp_{K(U, X|Z)} Z,$$

with  $K(E|Z)$  the uniform distribution on  $[0, 1]$ , and:

$$X = \tilde{R}(E|Z) \quad K(U, X|Z)\text{-a.s.},$$

for some measurable function  $\tilde{R}$ .

**Corollary 2.7.7.** *Let  $X$  and  $Z$  be random variables with values in any standard measurable spaces  $\mathcal{X}$  and  $\mathcal{Z}$ , resp., and with a joint distribution  $P(X, Z)$ . Then there exists a uniformly distributed random variable  $E$  on  $[0, 1]$  that is  $P$ -independent of  $Z$  and a measurable function  $g$  such that  $X = g(E, Z)$   $P$ -almost-surely.*

*Proof.* The regular conditional probability distribution  $P(X|Z)$  exists for standard measurable spaces (and is unique up to a  $P(Z)$ -zero-set), and is a Markov kernel. Then apply the result from above for  $K(X|Z) := P(X|Z)$  to get  $g(e, z) := \tilde{R}(e|z)$  and  $E$ .  $\square$

**Proofs - Deterministic Representation of Markov Kernels** In this section we generalize a few folklore results via now standard techniques that were introduced in [Dar53, Č82].

**Lemma 2.7.8.** *Let  $\bar{\mathbb{R}} := [-\infty, \infty]$  be endowed with the usual ordering and Borel  $\sigma$ -algebra. Let  $P$  be a probability measure on  $\bar{\mathbb{R}}$  and  $F(x) := P([-\infty, x])$ . Then  $F : \bar{\mathbb{R}} \rightarrow [0, 1]$  is non-decreasing, right-continuous with at most countably many discontinuities and  $F(\infty) = 1$ . So  $R(t) := \inf F^{-1}([t, 1])$  is a well-defined map  $R : [0, 1] \rightarrow \bar{\mathbb{R}}$ , non-decreasing, left-continuous with at most countably many discontinuities and  $R(0) = -\infty$ . Furthermore, for  $x \in \bar{\mathbb{R}}$  and  $t \in [0, 1]$  we have:*

$$t \leq F(x) \iff R(t) \leq x.$$

*In particular, we have  $F(R(t)) \geq t$ , thus  $R(t) \in F^{-1}([t, 1])$  the minimal element. We also have  $R(F(x)) \leq x$ , with equality if and only if  $x \in R([0, 1])$ . Furthermore,  $F$  and  $R$  are measurable and  $R_*\lambda = P$ . We also have that  $R$  is a reflexive generalized inverse of  $F$ , i.e.:*

$$F \circ R \circ F = F, \quad R \circ F \circ R = R.$$

*Proof.* From the properties of  $P$  it is clear that  $F$  is non-decreasing, right-continuous and  $F(\infty) = 1$ .

Let  $D_F \subseteq \bar{\mathbb{R}}$  be the set of discontinuities of  $F$  and  $x \in D_F$ . Then there exists a  $q(x) \in \mathbb{Q}$  such that  $F_-(x) < q(x) < F_+(x)$ . If now  $x_1 < x_2$  are two such points we get:

$q(x_1) < F_+(x_1) \leq F_-(x_2) < q(x_2)$ . So the map  $q : D_F \rightarrow \mathbb{Q}$  is injective. Thus  $D_F$  is countable.

Next, we show that  $R(t) \in F^{-1}([t, 1])$ , thus  $R(t) = \min F^{-1}([t, 1])$ . For this let  $(x_n)_{n \in \mathbb{N}} \subseteq F^{-1}([t, 1])$  be a non-increasing sequence converging to  $R(t)$ . Then by the right-continuity  $F(x_n)$  converges to  $F(R(t))$  from above. So we have:

$$F(R(t)) = \inf_{n \in \mathbb{N}} F(x_n) \geq t.$$

It follows that  $F(R(t)) \geq t$  and thus  $R(t) \in F^{-1}([t, 1])$ . This shows the claim.

$R$  is clearly non-decreasing, thus has only a countable set of discontinuities  $D_R \subseteq [0, 1]$  by the same arguments as before, and  $R(0) = -\infty$ . To see that  $R(t)$  is left-continuous let  $t \in [0, 1]$  and  $(t_n)_{n \in \mathbb{N}}$  a non-decreasing sequence converging to  $t$  from below. Then by the monotonicity of  $R$  we have  $\sup_{n \in \mathbb{N}} R(t_n) \leq R(t)$ . On the other hand we have:

$$t = \sup_{n \in \mathbb{N}} t_n \leq \sup_{n \in \mathbb{N}} F(R(t_n)) \leq F(\sup_{n \in \mathbb{N}} R(t_n)),$$

implying:  $\sup_{n \in \mathbb{N}} R(t_n) \in F^{-1}([t, 1])$  and thus  $\sup_{n \in \mathbb{N}} R(t_n) \geq R(t)$ , leading to equality, which shows the claim.

For any  $x \in \bar{\mathbb{R}}$  we have the implication:

$$x \geq R(t) \implies F(x) \geq F(R(t)) \geq t.$$

For any  $x \in \bar{\mathbb{R}}$  and any  $t \in [0, 1]$  we have the implications:

$$\begin{aligned} t \leq F(x) &\iff F(x) \in [t, 1] \\ &\iff x \in F^{-1}([t, 1]) \\ &\implies x \geq \inf F^{-1}([t, 1]) = R(t). \end{aligned}$$

Together this shows for any  $x \in \bar{\mathbb{R}}$  and  $t \in [0, 1]$  the equivalence:

$$t \leq F(x) \iff R(t) \leq x.$$

Since  $F(x) \leq F(x)$  we get  $R(F(x)) \leq x$  for all  $x \in \bar{\mathbb{R}}$ . If equality holds then  $x \in R([0, 1])$ . And, if  $x = R(t)$  for some  $t \in [0, 1]$  then we use the inequalities  $x \geq R(F(x))$  and  $F(R(t)) \geq t$  to conclude:

$$x \geq R(F(x)) = R(F(R(t))) \geq R(t) = x,$$

showing equality, and that:

$$R \circ F \circ R = R.$$

Similarly for  $t = F(x)$  we get:

$$t \leq F(R(t)) = F(R(F(x))) \leq F(x) = t,$$

showing

$$F \circ R \circ F = F.$$

Now consider the uniform distribution  $\lambda$  on  $[0, 1]$  and any  $x \in \bar{\mathbb{R}}$ . Then we have:

$$\begin{aligned}
(R_*\lambda)([-\infty, x]) &= \lambda(R^{-1}([-\infty, x])) \\
&= \lambda(t \in [0, 1] \mid R(t) \leq x) \\
&= \lambda(t \in [0, 1] \mid t \leq F(x)) \\
&= \lambda([0, F(x)]) \\
&= F(x) \\
&= P([-\infty, x]).
\end{aligned}$$

It follows that:  $R_*\lambda = P$ . □

**Lemma 2.7.9.** *Let the notation be like in lemma 2.7.8. For  $u \in [0, 1]$  and  $x \in \bar{\mathbb{R}}$  define:*

$$F_u(x) := E(x; u) := P([-\infty, x]) + u \cdot P(\{x\}).$$

*Then  $E : \bar{\mathbb{R}} \times [0, 1] \rightarrow [0, 1]$  is measurable, non-decreasing in both arguments with  $F_0(-\infty) = 0$ ,  $F_1(\infty) = 1$ ,  $F_0$  is left-continuous and*

$$F_u(\tilde{x}) \leq F_1(\tilde{x}) \leq F_0(x) \leq F_u(x)$$

*for any  $\tilde{x} < x$ ,  $u \in [0, 1]$ . We further have for every  $u \in (0, 1]$ :*

$$R \circ F_u \circ R = R,$$

*and  $R \circ F_u = \text{id}_{\bar{\mathbb{R}}}$   $P$ -almost-surely for any  $u \in (0, 1]$ .*

*Proof.* Most of the properties are clear from its definition. Let  $\tilde{x} < x$  then  $[-\infty, \tilde{x}] \subseteq [-\infty, x]$  and thus  $F_1(\tilde{x}) \leq F_0(x)$ .

To show  $R \circ F_u \circ R = R$  fix a  $t \in [0, 1]$ ,  $u \in (0, 1]$  and let  $x := R(t)$ . If  $F_1$  is continuous in  $x$  then  $F_u = F_1$  and the claim  $R \circ F_1 \circ R = R$  was already shown using the inequalities:

$$x \geq R(F_1(x)) = R(F_1(R(t))) \geq R(t) = x.$$

So let us assume that  $F_1$  is discontinuous in  $x = R(t)$ . Then  $F_u(x) \in (F_0(x), F_1(x)]$ . We have:

$$R(F_u(x)) = \min\{\tilde{x} \in \bar{\mathbb{R}} \mid F_1(\tilde{x}) \geq F_u(x)\}.$$

If  $F_1(\tilde{x}) \geq F_u(x) > F_0(x)$  then  $\tilde{x} \geq x$ , otherwise  $\tilde{x} < x$  leads to the contradiction  $F_1(\tilde{x}) \leq F_0(x)$ . Since clearly  $F_1(x) \geq F_u(x)$  we must have:

$$R(F_u(x)) = x,$$

with  $x = R(t)$ , which proves the claim:  $R \circ F_u \circ R = R$  for  $u \in (0, 1]$ .

We now want to show that  $R \circ F_u = \text{id}_{\bar{\mathbb{R}}}$   $P$ -a.s. for  $u \in (0, 1]$ . From  $R \circ F_u \circ R = R$  we already see that  $R \circ F_u|_{R([0,1])} = \text{id}_{R([0,1])}$ . We will see below that  $C := \bar{\mathbb{R}} \setminus R([0, 1])$  is measurable and  $P(C) = 0$ , which will prove the claim.

In the following we will only need  $F = F_1$ . First, by lemma 2.7.8 we know that for any

$x \in \bar{\mathbb{R}}$  we have  $R(F(x)) \leq x$  with equality if and only if  $x \in R([0, 1])$ . So this gives us the equivalence:

$$x \in C \iff x > R(F(x)).$$

We now claim that  $(R(F(x)), x] \subseteq C$  for every  $x \in C$ : Indeed, If  $\tilde{x} \in (R(F(x)), x]$  then:

$$F(x) = F(R(F(x))) \leq F(\tilde{x}) \leq F(x)$$

and thus  $F(\tilde{x}) = F(x)$ , from which follows that  $R(F(\tilde{x})) = R(F(x)) < \tilde{x}$  and ergo  $\tilde{x} \in C$ .

It follows that  $C$  is the union of such intervals  $(R(F(x)), x]$  with  $x \in C$ . Furthermore,  $F(C)$  is contained in the set of discontinuities  $D_R$  of  $R$ : otherwise there would be an  $x \in C$  and a  $t \geq F(x)$  such that  $R(t) \in (R(F(x)), x] \subseteq C$ , which is a contradiction. Since  $D_R$  is countable it must follow that  $F(C)$  and thus also  $R(F(C))$  is at most countable. Write  $R(F(C)) = \{x_n \mid n \in \mathbb{N}\}$ , which is the set of the possible left end-points of the above intervals. For each fixed  $n \in \mathbb{N}$  let

$$C_n := \{x \in C \mid R(F(x)) = x_n\},$$

which is, as a union of intervals  $(x_n, x]$ ,  $x \in C_n$ , either of the form  $(x_n, \bar{x}_n]$  or  $(x_n, \bar{x}_n)$  with  $\bar{x}_n := \sup C_n$ . In both cases we can cover  $C_n$  by  $C_{n,m} := (x_n, x_{n,m}]$  with  $x_{n,m} \in C_n$  either equal to  $\bar{x}_n$  or converging to it from below for running  $m$ . So we can write  $C$  as the countable union:

$$C = \bigcup_{n,m \in \mathbb{N}} C_{n,m}.$$

We now have for each  $x = x_{n,m}$ :

$$P(C_{n,m}) = P((x_n, x]) = P((R(F(x)), x]) = F(x) - F(R(F(x))) = F(x) - F(x) = 0.$$

This implies:

$$P(C) = P\left(\bigcup_{n,m \in \mathbb{N}} C_{n,m}\right) \leq \sum_{n,m \in \mathbb{N}} P(C_{n,m}) = 0,$$

showing that  $P(C) = 0$  and thus:

$$R \circ F_u = \text{id}_{\bar{\mathbb{R}}} \quad P\text{-a.s.}$$

for  $u \in (0, 1]$ . □

**Lemma 2.7.10.** *Let the notations be like in lemma 2.7.8 and lemma 2.7.9. Let  $\lambda$  be the uniform distribution on  $[0, 1]$  and  $\bar{P} := P \otimes \lambda$  the product distribution on  $\bar{\mathbb{R}} \times [0, 1]$ . For every  $e \in [0, 1]$  define the event:*

$$\{E \leq e\} := \{(x, u) \in \bar{\mathbb{R}} \times [0, 1] \mid E(x; u) \leq e\}.$$

Then  $\bar{P}(E \leq e) = e$ . In other words, the random variable:

$$\begin{aligned} E : \bar{\mathbb{R}} \times [0, 1] &\rightarrow [0, 1], \\ (x, u) &\mapsto P([-\infty, x]) + u \cdot P(\{x\}), \end{aligned}$$

is uniformly distributed under  $\bar{P} = P \otimes \lambda$ .

Furthermore,  $R(E) = X$   $\bar{P}$ -a.s., where  $X : \bar{\mathbb{R}} \times [0, 1] \rightarrow \bar{\mathbb{R}}$  is the canonical projection onto the first factor:  $X(x, u) := x$ , which has distribution  $P$ .

*Proof.* First, since  $\lambda(\{0\}) = 0$  we can w.l.o.g. exclude  $u = 0$  and restrict  $\bar{P}$  to  $\bar{\mathbb{R}} \times (0, 1]$ . We have seen in lemma 2.7.9 that  $R \circ F_u \circ R = R$  for  $u \in (0, 1]$ , which translates to:

$$R \circ E|_{R([0,1]) \times (0,1]} = X|_{R([0,1]) \times (0,1]}.$$

Also with  $C := \bar{\mathbb{R}} \setminus R([0, 1])$  we get:

$$\bar{P}(C \times (0, 1]) = P(C) \cdot \lambda((0, 1]) = 0 \cdot 1 = 0.$$

So we get the second claim that:

$$R \circ E = X \quad \bar{P}\text{-a.s.}$$

Now we turn to  $\{E \leq e\}$  for  $e \in [0, 1]$ . We abbreviate  $U : \bar{\mathbb{R}} \times [0, 1] \rightarrow [0, 1]$  to be the projection onto the second factor:  $U(x, u) := u$ , which is uniformly distributed under  $\bar{P}$ , and also  $p(x) := P(\{x\}) = F_1(x) - F_0(x)$ . With these notations:  $E = F_0(X) + U \cdot p(X)$ . First, we show that  $\bar{P}(E = e) = 0$  for all  $e \in [0, 1]$ . For this let  $x := R(e)$ . Then by the above ( $R(E) = X$   $\bar{P}$ -a.s.) we have:

$$\bar{P}(E = e) = \bar{P}(E = e, X = x).$$

We have to distinguish between two cases:  $p(x) = 0$  and  $p(x) > 0$ .

Case  $p(x) = 0$ : We have:

$$\begin{aligned} \bar{P}(E = e) &= \bar{P}(E = e, X = x) \\ &\leq \bar{P}(X = x) \\ &= p(x) \\ &= 0. \end{aligned}$$

Case  $p(x) > 0$ : We get:

$$\begin{aligned} \bar{P}(E = e) &= \bar{P}(E = e, X = x) \\ &= \bar{P}(F_0(X) + U \cdot p(X) = e, X = x) \\ &= \bar{P}\left(U = \frac{e - F_0(x)}{p(x)}, X = x\right) \\ &= \lambda\left(\left\{\frac{e - F_0(x)}{p(x)}\right\}\right) \cdot p(x) \\ &= 0. \end{aligned}$$

To prove  $\bar{P}(E \leq e) = e$  for  $e \in [0, 1]$  we have several cases:

Case  $e \in F_1(\bar{\mathbb{R}})$ : Let  $\tilde{x}$  be any element in  $\bar{\mathbb{R}}$  with  $e = F_1(\tilde{x})$  (e.g.  $\tilde{x} = R(e)$ ). Then we



get:

$$\begin{aligned}
\bar{P}(E \leq e) &= \bar{P}(E \leq F_1(\tilde{x})) \\
&= \bar{P}(R(E) \leq \tilde{x}) \\
&\stackrel{R \circ E = X}{=} \bar{P}(X \leq \tilde{x}) \\
&= P([-\infty, \tilde{x}] \cdot \lambda((0, 1])) \\
&= F_1(\tilde{x}) \cdot 1 \\
&= e.
\end{aligned}$$

For the cases  $e \notin F_1(\bar{\mathbb{R}})$  we put  $x := R(e)$  and  $\tilde{e} := F_0(x)$ .

Then by definition,  $x$  is minimal with  $F_1(x) \geq e$ . We also have  $\tilde{e} = F_0(x) \leq e$ . Otherwise:  $e < F_0(x) = \sup_{\tilde{x} < x} F_1(\tilde{x})$  implied that there existed  $\tilde{x} < x$  with  $e < F_1(\tilde{x}) \leq F_0(x)$ , which is a contradiction to the minimality of  $x = R(e)$ . Since  $\tilde{e} \leq e$  we can decompose:

$$\bar{P}(E \leq e) = \bar{P}(E < \tilde{e}) + \bar{P}(E = \tilde{e}) + \bar{P}(\tilde{e} < E \leq e).$$

We have already seen that the second term  $\bar{P}(E = \tilde{e}) = 0$  vanishes.

For the first term we have:

$$\begin{aligned}
\bar{P}(E < \tilde{e}) &= \bar{P}(E < F_0(x)) \\
&= \bar{P}(E < \sup_{\tilde{x} < x} F_1(\tilde{x})) \\
&= \sup_{\tilde{x} < x} \bar{P}(E \leq F_1(\tilde{x})) \\
&\stackrel{(*)}{=} \sup_{\tilde{x} < x} F_1(\tilde{x}) \\
&= F_0(x) \\
&= \tilde{e}.
\end{aligned}$$

Equation (\*) comes from the previous case for  $F_1(\tilde{x}) \in F_1(\bar{\mathbb{R}})$ .

For the third term  $\bar{P}(\tilde{e} < E \leq e)$  first note that  $E \in (\tilde{e}, e]$  implies that  $X = x$   $\bar{P}$ -a.s. by applying  $R$ : Indeed, every element  $t \in (\tilde{e}, e] \subseteq (F_0(x), F_1(x)]$  can be written as  $t = F_{\tilde{u}}(x)$  for an  $\tilde{u} \in (0, 1]$  and we can use:

$$R(t) = R(F_{\tilde{u}}(R(e))) = R(e) = x.$$

For  $p(x) > 0$  and the above we get:

$$\begin{aligned}
\bar{P}(\tilde{e} < E \leq e) &= \bar{P}(\tilde{e} < E \leq e, X = x) \\
&= \bar{P}(0 < F_0(X) + U \cdot p(X) - F_0(x) \leq e - \tilde{e}, X = x) \\
&= \bar{P}(0 < U \leq \frac{e - \tilde{e}}{p(x)}, X = x) \\
&= \lambda \left( \left( 0, \frac{e - \tilde{e}}{p(x)} \right] \right) \cdot P(\{x\}) \\
&= \frac{e - \tilde{e}}{p(x)} \cdot p(x) \\
&= e - \tilde{e}.
\end{aligned}$$

For the case  $p(x) = 0$ ,  $\bar{P}(\tilde{e} < E \leq e, X = x)$  can be upper bounded by  $\bar{P}(X = x) = p(x) = 0$  as before, but we also have  $\tilde{e} - e = 0$  in this case, and the equality stays trivially true as well.

Putting all together we get:

$$\begin{aligned}\bar{P}(E \leq e) &= \bar{P}(E < \tilde{e}) + \bar{P}(E = \tilde{e}) + \bar{P}(\tilde{e} < E \leq e) \\ &= \tilde{e} + 0 + e - \tilde{e} \\ &= e.\end{aligned}$$

This shows the claim. □

### 3. Graph Theory

#### 3.1. Core Concepts

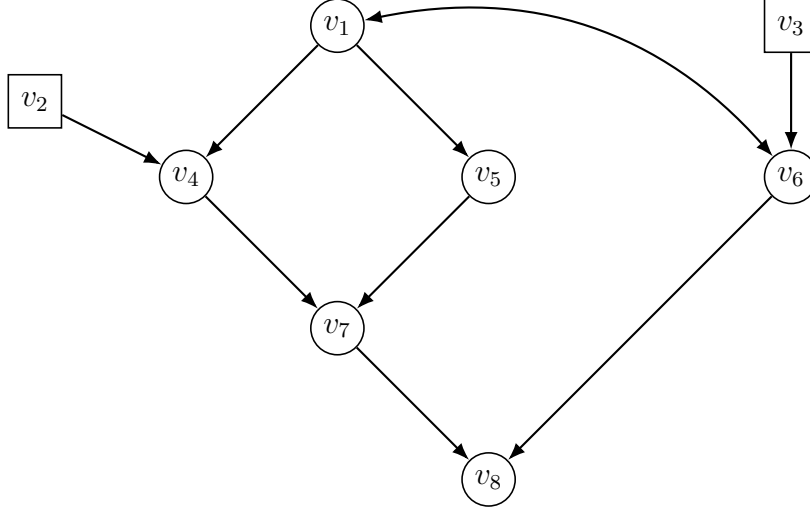


Figure 3: Conditional Acyclic Directed Mixed Graph (CADMG).

**Definition 3.1.1** (Conditional directed mixed graphs (CDMG)). A conditional directed mixed graph (CDMG)  $G$ —per definition—consists of two (disjoint) sets of vertices (also called nodes):

- i.)  $J$ , whose elements are called input nodes,
- ii.)  $V$ , whose elements are called output nodes,

and two (disjoint) sets of edges:

- iii.)  $E \subseteq (J \cup V) \times V$  the set of directed edges,
- iv.)  $L \subseteq V \times V / ((v_1, v_2) \sim (v_2, v_1))$ , the set of bi-directed edges,

$$\text{with: } (v_1, v_2) \in L \implies v_1 \neq v_2 \wedge (v_2, v_1) \in L.$$

**Notation 3.1.2.** Let  $G = (J, V, E, L)$  be a CDMG. We will write:

1.  $v \in G$  to mean  $v \in J \cup V$ ,
2.  $v_1 \rightarrow v_2 \in G$  to mean  $(v_1, v_2) \in E$ ,
3.  $v_1 \leftarrow v_2 \in G$  to mean  $(v_2, v_1) \in E$ ,
4.  $v_1 \leftrightarrow v_2 \in G$  to mean  $(v_1, v_2) \in L$ ,
5.  $v_1 * \rightarrow v_2 \in G$  to mean that either  $v_1 \rightarrow v_2 \in G$  or  $v_1 \leftrightarrow v_2 \in G$ ,

6.  $v_1 \leftarrow^* v_2 \in G$  to mean that either  $v_1 \leftarrow v_2 \in G$  or  $v_1 \leftrightarrow v_2 \in G$ ,
7.  $v_1 \rightarrow^* v_2 \in G$  to mean that either  $v_1 \rightarrow v_2 \in G$  or  $v_1 \leftarrow v_2 \in G$  or  $v_1 \leftrightarrow v_2 \in G$ .

The star stands for a placeholder to mean: “arrowhead or tail”.

**Definition 3.1.3.** Let  $G = (J, V, E, L)$  be a CDMG.

1. If  $v_1 \rightarrow^* v_2 \in G$  then we call  $v_1$  and  $v_2$  adjacent in  $G$ .
2. Edges of the form  $v_1 \leftarrow v_2$  or  $v_1 \leftrightarrow v_2$  are called into  $v_1$ .  
Edges of the form  $v_1 \rightarrow v_2$  or  $v_1 \leftrightarrow v_2$  are called into  $v_2$ .
3. Edges of the form  $v_1 \rightarrow v_2$  or  $v_2 \leftarrow v_1$  are called out of  $v_1$ .

**Remark 3.1.4.** With the notations 3.1.2 the restrictions in definition 3.1.1 mean that the nodes  $j \in J$  will not have any arrowheads pointing towards them:  $j \leftarrow^* v \notin G$ . Nodes  $j \in J$  can only point towards nodes  $v \in V$ : edges  $j \rightarrow v$  are allowed. Furthermore, no two nodes in  $J$  are adjacent.

**Definition 3.1.5 (Walks).** Let  $G = (J, V, E, L)$  be a CDMG and  $v, w \in G$ .

1. A walk from  $v$  to  $w$  in  $G$  is a finite alternating sequence of adjacent nodes and edges

$$v = v_0, a_0, v_1, \dots, v_{n-1}, a_{n-1}, v_n = w$$

in  $G$  for some  $n \geq 0$ , i.e. such that for every  $k = 0, \dots, n-1$  we have that  $a_k = (v_k, v_{k+1}) \in E \cup L$  or  $a_k = (v_{k+1}, v_k) \in E$ , and with end nodes  $v_0 = v$  and  $v_n = w$ . An example walk from  $v_0$  to  $v_3$  could look like:

$$v_0 \rightarrow v_1 \leftarrow v_2 \leftrightarrow v_3, \quad \text{with} \quad v_0 \rightarrow v_1, v_2 \rightarrow v_1 \in E, v_2 \leftrightarrow v_3 \in L.$$

The same node may appear multiple times in a walk. Also the trivial walk consisting of a single node  $v_0 \in G$  is allowed (if  $v = w$ ). The walk is called into  $v_0$  if  $a_0 = v_0 \leftarrow^* v_1$ , and out of  $v_0$  if  $a_0 = v_0 \rightarrow v_1$ . Similarly, it is called into  $v_n$  if  $a_{n-1} = v_{n-1} \rightarrow^* v_n$  and out of  $v_n$  if  $a_{n-1} = v_{n-1} \leftarrow v_n$ .

2. A directed walk from  $v$  to  $w$  in  $G$  is of the form:

$$v = v_0 \rightarrow v_1 \rightarrow \dots \rightarrow v_{n-1} \rightarrow v_n = w,$$

for some  $n \geq 0$ , where all arrowheads point in the direction of  $w$  and there are no arrowheads pointing back.

3. A bi-directed walk from  $v$  to  $w$  in  $G$  is of the form:

$$v = v_0 \leftrightarrow v_1 \leftrightarrow \dots \leftrightarrow v_{n-1} \leftrightarrow v_n = w,$$

for some  $n \geq 0$ , where all edges are bi-directed.

4. A collider walk from  $v$  to  $w$  in  $G$  is of the form:

$$v = v_0 \ast \rightarrow v_1 \leftrightarrow \dots \leftrightarrow v_{n-1} \leftarrow \ast v_n = w,$$

for some  $n \geq 0$ , where all nodes in between  $v$  and  $w$  have two arrowheads pointing towards them (a.k.a. collider). Note that for  $n = 1$  this reads:  $v \ast \ast w \in G$ .

5. A walk is called path if no node occurs more than once.

6. A bifurcation from  $v$  to  $w$  in  $G$  is a walk of the form:

$$v = v_0 \leftarrow v_1 \leftarrow \dots \leftarrow v_{k-1} \leftarrow \ast v_k \rightarrow \dots \rightarrow v_{n-1} \rightarrow v_n = w,$$

such that  $v \neq w$ , the walk contains both endnodes exactly once, every node has at most one arrowhead pointing towards it, and both endnodes have exactly one arrowhead pointing towards them. If the edge  $v_{k-1} \leftarrow \ast v_k$  is directed ( $v_{k-1} \leftarrow v_k$ ) then we say that the bifurcation has source  $v_k$ .

**Definition 3.1.6** (Family relationships). Let  $G = (J, V, E, L)$  be a CDMG,  $v, w \in V$  and  $A \subseteq J \cup V$  a subset of nodes. We then define:

1. The set of parents of  $v$  in  $G$ :

$$\text{Pa}^G(v) := \{w \in G \mid w \rightarrow v \in G\}.$$

The set of parents of  $A$  in  $G$ :

$$\text{Pa}^G(A) := \bigcup_{v \in A} \text{Pa}^G(v).$$

2. The set of children of  $v$  in  $G$ :

$$\text{Ch}^G(v) := \{w \in G \mid v \rightarrow w \in G\}.$$

The set of children of  $A$  in  $G$ :

$$\text{Ch}^G(A) := \bigcup_{v \in A} \text{Ch}^G(v).$$

3. The set of siblings of  $v$  in  $G$ :

$$\text{Sib}^G(v) := \{w \in G \mid v \leftrightarrow w \in G\}.$$

4. The set of ancestors of  $v$  in  $G$ :

$$\text{Anc}^G(v) := \{w \in G \mid \exists \text{ directed walk: } w \rightarrow \dots \rightarrow v \in G\}.$$

Note:  $v \in \text{Anc}^G(v)$ .

The set of ancestors of  $A$  in  $G$ :

$$\text{Anc}^G(A) := \bigcup_{v \in A} \text{Anc}^G(v).$$

Note:  $A \subseteq \text{Anc}^G(A)$ .

5. The set of descendants of  $v$  in  $G$ :

$$\text{Desc}^G(v) := \{w \in G \mid \exists \text{ directed walk: } v \rightarrow \dots \rightarrow w \in G\}.$$

Note:  $v \in \text{Desc}^G(v)$ .

The set of descendants of  $A$  in  $G$ :

$$\text{Desc}^G(A) := \bigcup_{v \in A} \text{Desc}^G(v).$$

Note:  $A \subseteq \text{Desc}^G(A)$ .

6. The set of non-descendants of  $A$  in  $G$ :

$$\text{NonDesc}^G(A) := (J \cup V) \setminus \text{Desc}^G(A).$$

7. The strongly connected component of  $v$  in  $G$ :

$$\text{Sc}^G(v) := \text{Anc}^G(v) \cap \text{Desc}^G(v).$$

Note:  $v \in \text{Sc}^G(v)$ .

The (union of) strongly connected components of  $A$  in  $G$ :

$$\text{Sc}^G(A) := \bigcup_{v \in A} \text{Sc}^G(v).$$

Note:  $A \subseteq \text{Sc}^G(A)$ .

8. The district of  $v$  in  $G$ :

$$\text{Dist}^G(v) := \{w \in G \mid \exists \text{ bi-directed walk: } v \leftrightarrow v_1 \leftrightarrow \dots \leftrightarrow v_{n-1} \leftrightarrow w \in G\}.$$

Note:  $v \in \text{Dist}^G(v)$ .

The district of  $A$  in  $G$ :

$$\text{Dist}^G(A) := \bigcup_{v \in A} \text{Dist}^G(v).$$

Note:  $A \subseteq \text{Dist}^G(A)$ .

**Definition 3.1.7** (Acyclicity). A CDMG  $G = (J, V, E, L)$  is called acyclic if there does not exist any non-trivial directed walk from  $v$  to itself in  $G$  for any node  $v \in G$ .

**Definition 3.1.8.** A Conditional Directed Mixed Graph (CDMG)  $G = (J, V, E, L)$  is called:

1. Conditional Acyclic Directed Mixed Graph (CADMG) if  $G$  is acyclic.
2. Directed Mixed Graph (DMG) if  $J = \emptyset$ .
3. Acyclic Directed Mixed Graph (ADMG) if  $G$  is acyclic and  $J = \emptyset$ .
4. Conditional Directed Graph (CDG) if  $L = \emptyset$ .
5. Directed Graph (DG) if  $J = \emptyset$  and  $L = \emptyset$ .
6. Conditional Directed Acyclic Graph (CDAG) if  $G$  is acyclic and  $L = \emptyset$ .
7. Directed Acyclic Graph (DAG) if  $G$  is acyclic,  $J = \emptyset$  and  $L = \emptyset$ .

**Definition 3.1.9** (Topological order). Let  $G = (J, V, E, L)$  be a CDMG. A topological order of  $G$  is a total order  $<$  of  $J \cup V$  such that for all  $v, w \in G$ :

$$v \in \text{Pa}^G(w) \implies v < w.$$

Equivalently, it can be described as an indexing of the nodes  $J \cup V = \{v_1, \dots, v_K\}$  where parents always precede their children.

**Lemma 3.1.10.** A CDMG  $G = (J, V, E, L)$  is acyclic if and only if it has a topological order.

**Definition 3.1.11** (Predecessors). Let  $G = (J, V, E, L)$  be a CDMG and  $<$  a total order of  $J \cup V$ . The set of predecessors of  $v$  in  $G$  are:

$$\text{Pred}_{<}^G(v) := \{w \in G \mid w < v\}.$$

We also put:

$$\text{Pred}_{\leq}^G(v) := \{w \in G \mid w < v\} \cup \{v\}.$$

## 3.2. Operations on Graphs

### 3.2.1. Hard Interventions on Graphs

**Definition 3.2.1** (Hard intervention on CDMGs). Let  $G = (J, V, E, L)$  be a CDMG and  $W \subseteq J \cup V$  a subset of nodes. The intervened CDMG w.r.t.  $W$  of  $G$  is the CDMG:

$$G_{\text{do}(W)} := (J_{\text{do}(W)}, V_{\text{do}(W)}, E_{\text{do}(W)}, L_{\text{do}(W)}),$$

where:

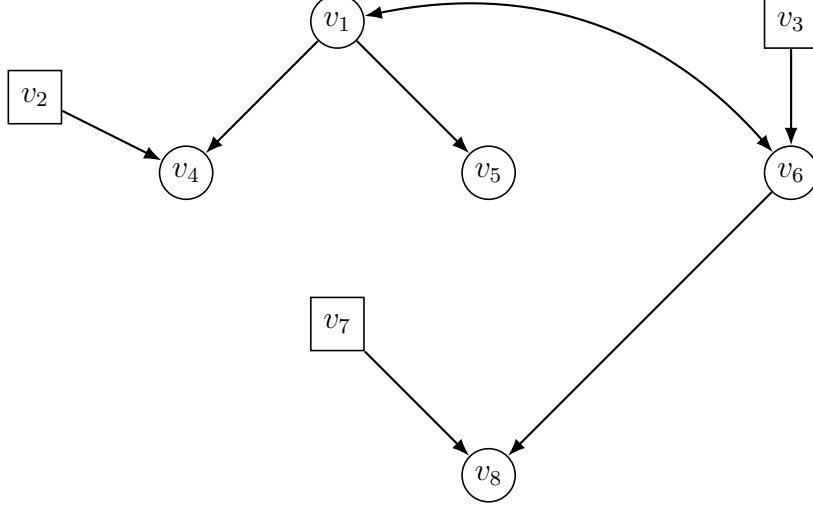


Figure 4: The CADMG from Figure 3 after hard intervention on node  $v_7$ .

- i.)  $J_{\text{do}(W)} := J \cup W$ ,
- ii.)  $V_{\text{do}(W)} := V \setminus W$ ,
- iii.)  $E_{\text{do}(W)} := E \setminus \{v \rightarrow w \mid v \in G, w \in W\}$ ,
- iv.)  $L_{\text{do}(W)} := L \setminus \{v \leftrightarrow w \mid v \in G, w \in W\}$ ,

where we turn all nodes from  $W$  into input nodes and remove all edges into nodes from  $W$ .

**Remark 3.2.2.** If  $G$  is acyclic then also  $G_{\text{do}(W)}$  is acyclic and a topological order for  $G$  is also one for  $G_{\text{do}(W)}$ .

**Lemma 3.2.3** (Hard interventions commute). Let  $G := (J, V, E, L)$  be a CDMG and  $W_1, W_2 \subseteq J \cup V$  two subsets of nodes from  $G$ . Then we have:

$$(G_{\text{do}(W_1)})_{\text{do}(W_2)} = (G_{\text{do}(W_2)})_{\text{do}(W_1)} = G_{\text{do}(W_1 \cup W_2)}.$$

The following proposition expresses the existence of a bifurcation with a source in terms of ancestral relations in intervened graphs.

**Proposition 3.2.4.** Let  $G = (J, V, E, L)$  be a CMDG. For  $v, w, c \in V \cup J$ : there exists a bifurcation between  $v$  and  $w$  in  $G$  with source  $c$  if and only if  $v \neq w$  and  $c \in \text{Anc}^{G_{\text{do}(w)}}(v) \setminus \{v\}$  and  $c \in \text{Anc}^{G_{\text{do}(v)}}(w) \setminus \{w\}$ .

*Proof.* A bifurcation between  $v$  and  $w$  with source  $c$  is a walk in  $G$  of the form  $v \leftarrow \dots \leftarrow c \rightarrow \dots \rightarrow w$ , where both  $v$  and  $w$  appear exactly once on the walk. This shows “ $\implies$ ”. For the other implication, note that  $c \in \text{Anc}^{G_{\text{do}(w)}}(v) \setminus \{v\}$  implies that there is a non-trivial directed path from  $c$  to  $v$  that does not pass through  $w$ . Similarly,  $c \in \text{Anc}^{G_{\text{do}(v)}}(w) \setminus \{w\}$  implies that there is a non-trivial directed path from  $c$  to  $w$  that does not pass through  $v$ . The concatenation of the two paths  $v \leftarrow \dots \leftarrow c \rightarrow \dots \rightarrow w$  is then a bifurcation between  $v$  and  $w$  with source  $c$ .  $\square$



### 3.2.2. Node Splitting Interventions on Graphs

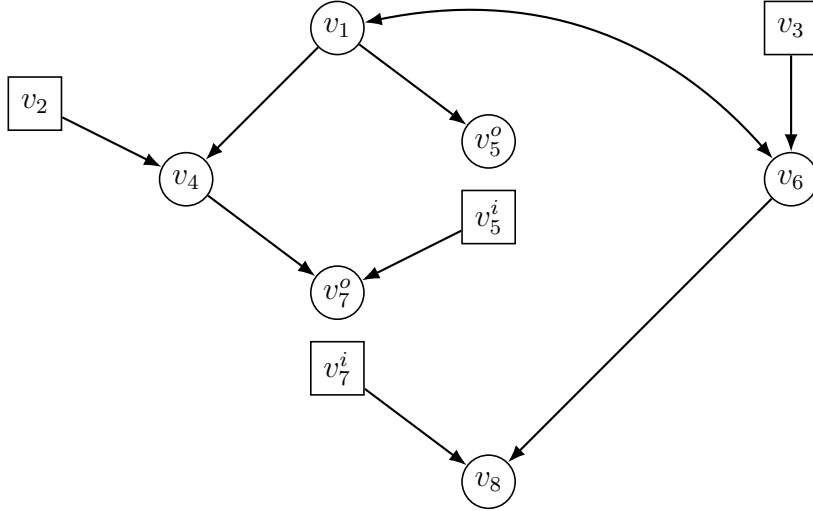


Figure 5: The CADMG from Figure 3 after a node-splitting hard intervention on  $v_5$  and  $v_7$ .

In this subsection we introduce *node-splitting hard interventions*. They were introduced for the purpose of representing *single-world intervention graphs (SWIGs)*, which represent the same output variable both before and after a hard intervention, see [RR13a, RR13b].

**Definition 3.2.5** (Node-splitting hard intervention on CADMGs). *Let  $G = (J, V, E, L)$  be a CADMG and  $W \subseteq V$  a subset of the output nodes. The single-world intervention graph (SWIG) w.r.t.  $W$  of  $G$  is the CADMG:*

$$G_{\text{swig}(W)} := (J_{\text{swig}(W)}, V_{\text{swig}(W)}, E_{\text{swig}(W)}, L_{\text{swig}(W)}),$$

constructed as follows. We first make two disjoint copies of the nodes in  $W$ :

$$W^o := \{w^o \mid w \in W\}, \quad W^i := \{w^i \mid w \in W\}.$$

Note that we consider  $w^o \neq w^i$  for  $w \in W$ . However, for brevity, for  $v \in J \cup V \setminus W$  we put:

$$v^o := v^i := v.$$

We then define:

- i.)  $J_{\text{swig}(W)} := J \dot{\cup} W^i,$
- ii.)  $V_{\text{swig}(W)} := (V \setminus W) \dot{\cup} W^o,$
- iii.)  $E_{\text{swig}(W)} := \{v_1^i \rightarrow v_2^o \mid v_1 \rightarrow v_2 \in E\},$
- iv.)  $L_{\text{swig}(W)} := \{v_1^o \leftrightarrow v_2^o \mid v_1 \leftrightarrow v_2 \in L\}.$

where we turn all nodes of  $W^i$  into input nodes, removing all edges into  $W^i$ , and we turn all nodes of  $W^o$  into output nodes, removing all edges out of  $W^o$ .

**Remark 3.2.6.** For a CADMG  $G = (J, V, E, L)$ , also  $G_{\text{swig}(W)}$  is acyclic. If  $<$  is any topological order of  $G$  given by enumerating all nodes  $v \in J \cup V$  via:

$$v_1 < v_2 < \dots < v_n,$$

then, for instance, a topological order for  $G_{\text{swig}(W)}$  can be achieved by assigning for a node  $v_j \in W$  with index  $j$  the node  $v_j^o$  the index  $j - \frac{1}{3}$  and  $v_j^i$  the index  $j + \frac{1}{3}$ , and then ordering all nodes according to their index value.

**Lemma 3.2.7** (Two disjoint node-splitting hard interventions commute). *Let  $G = (J, V, E, L)$  be a CADMG and  $W_1, W_2 \subseteq V$  two disjoint subsets of the output nodes from  $G$ . Then the CADMG obtained from first node-splitting on  $W_1$  and then node-splitting on  $W_2$  is the same CADMG that arises from first node-splitting on  $W_2$  and then node-splitting on  $W_1$ :*

$$(G_{\text{swig}(W_1)})_{\text{swig}(W_2)} = (G_{\text{swig}(W_2)})_{\text{swig}(W_1)} = G_{\text{swig}(W_1 \cup W_2)}.$$

**Lemma 3.2.8** (Disjoint hard interventions and node-splitting hard interventions commute). *Let  $G = (J, V, E, L)$  be a CADMG and  $W_1 \subseteq J \cup V$  and  $W_2 \subseteq V$  two disjoint subsets of nodes from  $G$ . Then the CADMG obtained from first hard intervening on  $W_1$  and then node-splitting on  $W_2$  is the same CADMG that arises from first node-splitting on  $W_2$  and then hard intervening on  $W_1$ .*

$$(G_{\text{do}(W_1)})_{\text{swig}(W_2)} = (G_{\text{swig}(W_2)})_{\text{do}(W_1)}.$$

**Remark 3.2.9.** Note that if  $W_1$  and  $W_2$  are not disjoint and  $w \in W_1 \cap W_2 \subseteq V$  then first hard intervening on  $w$  turns  $w$  into an input node, for now indicated as  $w^i$ , and a node-splitting hard intervention (if we would define it for input nodes) would not change  $w^i$ . If, on the other hand, we would first split the node  $w$  into  $w^o$  and  $w^i$  then we would first need to resolve the ambiguity on which of those two the hard intervention should be applied. A hard intervention on  $w^i$  would not do anything, but would leave the additional output node  $w^o$  in the graph, while hard intervening on  $w^o$  would turn  $w^o$  into an additional input node, for now indicated as  $(w^o)^i$ . So in the latter case we are left with two input node  $(w^o)^i$ , which does not have any edges, and  $w^i$ , which might have outgoing edges.

### 3.2.3. Intervention Nodes

More generally, interventions (both hard and soft) can be modeled graphically via auxiliary intervention nodes.

**Definition 3.2.10** (Extending CDMGs with intervention nodes). *Let  $G = (J, V, E, L)$  be a CDMG and  $W \subseteq J \cup V$  a subset of nodes. The extended CDMG of  $G$  w.r.t. nodes  $W \subseteq J \cup V$  and corresponding intervention nodes  $I_W = \{I_w \mid w \in W\}$  is the CDMG:*

$$G_{\text{do}(I_W)} := (J_{\text{do}(I_W)}, V_{\text{do}(I_W)}, E_{\text{do}(I_W)}, L_{\text{do}(I_W)}),$$

where:

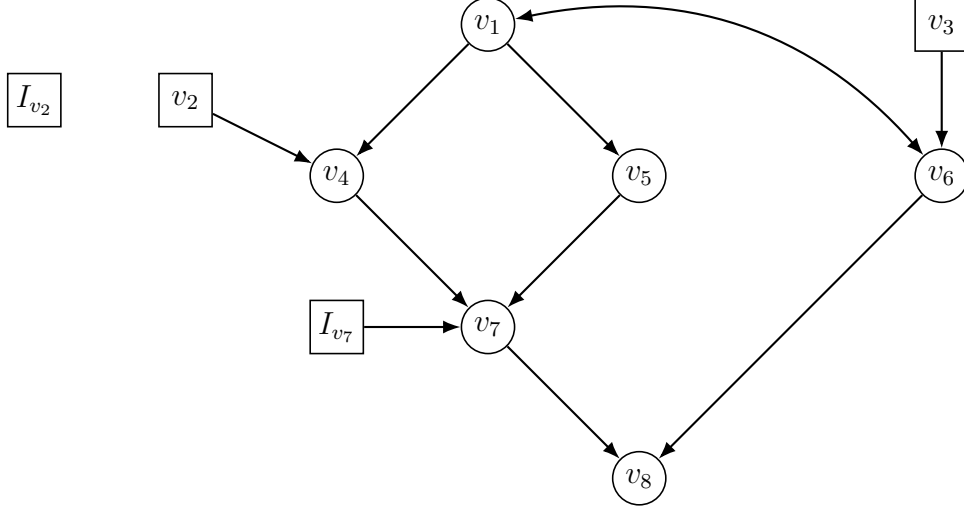


Figure 6: The CADMG from Figure 3 with additional intervention node  $I_{v_7}$  at node  $v_7$ .

- i.)  $J_{\text{do}(I_W)} := J \dot{\cup} \{I_w \mid w \in W\}$ ,
- ii.)  $V_{\text{do}(I_W)} := V$ ,
- iii.)  $E_{\text{do}(I_W)} := E \dot{\cup} \{I_w \rightarrow w \mid w \in W \setminus J\}$ ,
- iv.)  $L_{\text{do}(I_W)} := L$ ,

where we just add edges  $I_w \rightarrow w$  for  $w \in W \setminus J$ , where  $I_w$  represent intervention nodes.

**Remark 3.2.11.** If a CDMG  $G = (J, V, E, L)$  is acyclic then also  $G_{\text{do}(I_W)}$  is acyclic and a topological order for  $G_{\text{do}(I_W)}$  is also one for  $G$ . Any topological order of  $G$  can be extended to one for  $G_{\text{do}(I_W)}$ , e.g. by putting all the  $I_w$  nodes first in the ordering.

**Lemma 3.2.12** (Adding intervention nodes commutes with disjoint hard interventions). Let  $G = (J, V, E, L)$  be a CDMG and  $W_1, W_2 \subseteq J \cup V$  two disjoint subsets of nodes from  $G$ . Then we have:

$$\left(G_{\text{do}(I_{W_1})}\right)_{\text{do}(I_{W_2})} = \left(G_{\text{do}(I_{W_2})}\right)_{\text{do}(I_{W_1})} = G_{\text{do}(I_{W_1 \cup W_2})}.$$

We also have:

$$\left(G_{\text{do}(I_{W_1})}\right)_{\text{do}(W_2)} = \left(G_{\text{do}(W_2)}\right)_{\text{do}(I_{W_1})} = G_{\text{do}(I_{W_1}, W_2)}.$$

**Lemma 3.2.13** (Adding intervention nodes commutes with disjoint node-splitting hard interventions). Let  $G = (J, V, E, L)$  be a CADMG and  $W_1 \subseteq V$  and  $W_2 \subseteq J \cup V$  two disjoint subsets of nodes from  $G$ . Then the CADMG that arises from first introducing intervention nodes  $I_{W_2}$  and then splitting the nodes from  $W_1$  is the same as the CADMG that arises from first splitting the nodes from  $W_1$  and then introducing the intervention nodes  $I_{W_2}$ :

$$\left(G_{\text{swig}(W_1)}\right)_{\text{do}(I_{W_2})} = \left(G_{\text{do}(I_{W_2})}\right)_{\text{swig}(W_1)}.$$

### 3.2.4. Marginalization of Graphs

**Definition 3.2.14** (Marginalization a.k.a. latent projection on CDMGs). *Let  $G = (J, V, E, L)$  be a CDMG and  $W \subseteq V$  a subset of output nodes. Then the marginalization of  $G$  w.r.t.  $W$  or the latent projection of  $G$  onto  $J \cup V \setminus W$  is the CDMG:*

$$G^{V \setminus W} := G^{\setminus W} := (J^{\setminus W}, V^{\setminus W}, E^{\setminus W}, L^{\setminus W}),$$

where:

- i.)  $J^{\setminus W} := J$ ,
- ii.)  $V^{\setminus W} := V \setminus W$ ,
- iii.)  $E^{\setminus W}$  consists of all directed edges  $\underline{v} \rightarrow \bar{v}$  with  $\underline{v}, \bar{v} \in J \cup V \setminus W$  for which there exists a directed walk in  $G$ :

$$\underline{v} \rightarrow w_1 \rightarrow \cdots \rightarrow w_{n-1} \rightarrow \bar{v},$$

where all intermediate nodes  $w_1, \dots, w_{n-1} \in W$  (if any).<sup>12</sup>

- iv.)  $L^{\setminus W}$  consists of all bi-directed edges  $\underline{v} \leftrightarrow \bar{v}$  with  $\underline{v}, \bar{v} \in V \setminus W$ ,  $\underline{v} \neq \bar{v}$ , for which there exists a bifurcation in  $G$ :

$$\underline{v} \leftarrow w_1 \leftarrow \cdots \leftarrow w_{k-1} \leftarrow^* w_k \rightarrow \cdots \rightarrow w_{n-1} \rightarrow \bar{v},$$

where all intermediate nodes  $w_1, \dots, w_{n-1} \in W$  (if any).

**Remark 3.2.15.** *Marginalization preserves ancestral relations, bifurcations and acyclicity:*

1. For  $v_1, v_2 \in G$  with  $v_1, v_2 \notin W$  we have the equivalence:

$$v_1 \in \text{Anc}^G(v_2) \iff v_1 \in \text{Anc}^{G^{\setminus W}}(v_2).$$

2. For  $v_1, v_2 \in G$  with  $v_1, v_2 \notin W$  there is a bifurcation between  $a$  and  $b$  in  $G$  if and only if there is a bifurcation between  $a$  and  $b$  in  $G^{\setminus W}$ .
3. If the CDMG  $G$  is acyclic then so is  $G^{\setminus W}$  and a topological order of  $G$  induces a topological order on  $G^{\setminus W}$  (by just ignoring the nodes from  $W$ ).

*Proof.* We prove 2. Let

$$v_1 = w_0 \leftarrow w_1 \leftarrow \cdots \leftarrow w_{k-1} \leftarrow^* w_k \rightarrow \cdots \rightarrow w_{n-1} \rightarrow w_n = v_2,$$

be a bifurcation in  $G$ . If one marginalizes out a single node  $u$  that is not on the bifurcation, then the same bifurcation exists in  $G^{\setminus \{u\}}$ . If one marginalizes out a single node  $u$  that appears on the bifurcation (but that is not an endpoint) one obtains again a bifurcation in  $G^{\setminus \{u\}}$ . The statement follows by induction.  $\square$

<sup>12</sup>Note that this may introduce self-cycles.

**Lemma 3.2.16** (Marginalizations commute). *Let  $G = (J, V, E, L)$  be a CDMG and  $W_1, W_2 \subseteq V$  two disjoint subsets of output nodes. Then we have:*

$$(G \setminus W_1) \setminus W_2 = (G \setminus W_2) \setminus W_1 = G \setminus (W_1 \cup W_2).$$

**Lemma 3.2.17** (Marginalization and intervention commute). *Let  $G = (J, V, E, L)$  be a CDMG and  $W_1 \subseteq J \cup V$  and  $W_2 \subseteq V$  two disjoint subsets of nodes from  $G$ . Then we have:*

$$(G_{\text{do}(W_1)}) \setminus W_2 = (G \setminus W_2)_{\text{do}(W_1)}.$$

*A similar statement holds for marginalizations and adding intervention nodes, and also for marginalizations and node-splitting interventions.*

**Lemma 3.2.18** (Marginalizing out the output part of splitted nodes equals hard intervention). *Let  $G = (J, V, E, L)$  be a CDMG and  $W \subseteq V$  be a subset of output nodes from  $G$ . Then the CDMG that arises by first splitting the nodes on  $W$  and then marginalizing out the nodes from  $W^o$  can be identified with the CDMG that arises by hard intervention on  $W$ :*

$$G_{\text{do}(W)} \cong (G_{\text{swig}(W)}) \setminus W^o, \quad w \mapsto w^i.$$

### 3.3. $\sigma$ -Separation

**Definition 3.3.1** (Colliders and non-colliders). *Let  $G = (J, V, E, L)$  be a CDMG and  $\pi$  a walk in  $G$ :*

$$\pi = (v_0 \xrightarrow{*} \dots \xrightarrow{*} v_n).$$

*A node  $v_k$ , or more precisely, the position  $k \in \{0, \dots, n\}$ , on the walk  $\pi$  is called:*

1. *a non-collider on  $\pi$ , if there is at most one arrowhead pointing towards  $v_k$ , i.e. if it falls into one of the following cases:*

$$\begin{array}{ll} \text{end-node:} & k \in \{0, n\}, \\ \text{left chain:} & v_{k-1} \leftarrow v_k \xleftarrow{*} v_{k+1}, \\ \text{right chain:} & v_{k-1} \xrightarrow{*} v_k \rightarrow v_{k+1}, \\ \text{fork:} & v_{k-1} \leftarrow v_k \rightarrow v_{k+1}; \end{array}$$

2. *a collider on  $\pi$ , if it is of the form:*

$$v_{k-1} \xrightarrow{*} v_k \xleftarrow{*} v_{k+1},$$

*i.e. if there are two arrowheads pointing towards  $v_k$  on the walk  $\pi$ .*

**Definition 3.3.2** (Blockable and unblockable non-colliders). *Let  $G = (J, V, E, L)$  be a CDMG and  $\pi$  a walk in  $G$ :*

$$\pi = (v_0 \xrightarrow{*} \dots \xrightarrow{*} v_n).$$

We call a non-collider  $v_k$  on  $\pi$  an *unblockable non-collider on  $\pi$*  if it is not an end-node ( $k \notin \{0, n\}$ ) and it only has outgoing edges on  $\pi$  to nodes in the same strongly connected component of  $G$ . That is, it is one of the following patterns:

$$\begin{aligned} \text{left chain:} \quad & v_{k-1} \leftarrow v_k \leftarrow^* v_{k+1} \quad \text{with} \quad v_{k-1} \in \text{Sc}^G(v_k) \\ \text{right chain:} \quad & v_{k-1} \xrightarrow{*} v_k \rightarrow v_{k+1} \quad \text{with} \quad v_{k+1} \in \text{Sc}^G(v_k) \\ \text{fork:} \quad & v_{k-1} \leftarrow v_k \rightarrow v_{k+1} \quad \text{with} \quad v_{k-1} \in \text{Sc}^G(v_k) \wedge v_{k+1} \in \text{Sc}^G(v_k) \end{aligned}$$

Otherwise,  $v_k$  is called a *blockable non-collider on  $\pi$* . This means that  $v_k$  is either an end-node ( $k \in \{0, n\}$ ) or it has at least one outgoing arrow  $v_k \rightarrow v_{k\pm 1}$  pointing to a node  $v_{k\pm 1}$  that lies in a different strongly connected component than  $v_k$ , i.e.  $v_{k\pm 1} \notin \text{Sc}^G(v_k)$ .

**Remark 3.3.3.** If  $G$  is acyclic then all non-colliders are blockable.

**Definition 3.3.4** ( $\sigma$ -blocked walks). Let  $G = (J, V, E, L)$  be a CDMG and  $C \subseteq J \cup V$  a subset of nodes and  $\pi$  a walk in  $G$ :

$$\pi = (v_0 \xrightarrow{*} \dots \xrightarrow{*} v_n).$$

We say that the walk  $\pi$  is:

1.  $C$ - $\sigma$ -open (or  $\sigma$ -open given  $C$ ) if and only if:
  - i.) all colliders  $v_k$  on  $\pi$  are in  $\text{Anc}^G(C)$ , and:
  - ii.) all blockable non-colliders  $v_k$  on  $\pi$  are not in  $C$ .
2.  $C$ - $\sigma$ -blocked (or  $\sigma$ -blocked given  $C$ ) if and only if:
  - i.) there exists a collider  $v_k$  on  $\pi$  that is not in  $\text{Anc}^G(C)$ , or:
  - ii.) there exists a blockable non-collider  $v_k$  on  $\pi$  in  $C$ .

Note that unblockable non-colliders are always  $C$ - $\sigma$ -open, regardless of the subset  $C \subseteq V \cup J$ .

**Definition 3.3.5** ( $\sigma$ -separation). Let  $G = (J, V, E, L)$  be a CDMG and  $A, B, C \subseteq J \cup V$  (not necessarily disjoint) subset of nodes. We then say that:

1.  $A$  is  $\sigma$ -separated from  $B$  given  $C$  in  $G$ , in symbols:

$$A \underset{G}{\perp}^{\sigma} B \mid C,$$

if every walk from a node in  $A$  to a node in  $J \cup B$  (sic!)<sup>13</sup> is  $C$ - $\sigma$ -blocked by  $C$ .

---

<sup>13</sup>The choice to include  $J$  here in this place is non-standard in the literature. However, if we include  $J$  in this definition here the implied (asymmetric) separoid rules for d-/ $\sigma$ -separation will be of the same form as those for Markov kernels regarding conditional independence. This is the reason we include  $J$  here.

2. If that property does not hold we will write:

$$A \not\perp_G^\sigma B | C.$$

3. We also define the special case:

$$A \perp_G^\sigma B \quad : \iff \quad A \perp_G^\sigma B | \emptyset.$$

The following result is often helpful to simplify proofs and making checking  $\sigma$ -separation feasible in practice.

**Proposition 3.3.6.** *Let  $G = (J, V, E, L)$  be a CDMG. For  $C \subseteq J \cup V$ , and  $w_1, w_2 \in J \cup V$ , the following are equivalent:*

1. *there exists a  $C$ - $\sigma$ -open path between  $w_1$  and  $w_2$  in  $G$ ;*
2. *there exists a  $C$ - $\sigma$ -open walk between  $w_1$  and  $w_2$  in  $G$ ;*
3. *there exists a  $C$ - $\sigma$ -open walk between  $w_1$  and  $w_2$  in  $G$  such that all its colliders lie in  $C$  (and not just in  $\text{Anc}^G(C)$ ).*

**Remark 3.3.7.** 1. *By Proposition 3.3.6 we have that  $A \perp_G^\sigma B | C$  is equivalent to either of the following:*

- a) *every walk from a node in  $A$  to a node in  $J \cup B$  is  $C$ - $\sigma$ -blocked by  $C$ ;*
- b) *every path from a node in  $A$  to a node in  $J \cup B$  is  $C$ - $\sigma$ -blocked by  $C$ .*

2. *Proposition 3.3.6 also shows that if  $A \not\perp_G^\sigma B | C$  holds then:*

- a) *there exists a (shortest)  $C$ - $\sigma$ -open path from a node in  $A$  to a node in  $J \cup B$ ;*
- b) *there exists a (shortest)  $C$ - $\sigma$ -open walk from a node in  $A$  to a node in  $J \cup B$  such that all its colliders lie in  $C$ .*

*In practice we usually check if every path is  $C$ - $\sigma$ -blocked or not. This is because there are, in contrast to walks, only a finite number of paths in a (finite) graph. In proofs, though, it often is easier to make use of walks, since these can be concatenated into walks (while one cannot in general concatenate two paths and again obtain a path).*

**Lemma 3.3.8** ( $\sigma$ -separation under marginalization). *Let  $G = (J, V, E, L)$  be a CDMG,  $A, B, C \subseteq J \cup V$  and  $D \subseteq V$  be subsets of nodes such that:*

$$D \cap (A \cup B \cup C) = \emptyset.$$

*Then we have the equivalence:*

$$A \perp_G^\sigma B | C \quad \iff \quad A \perp_{G \setminus D}^\sigma B | C.$$

**Remark 3.3.9.** *If a CDMG  $G$  is acyclic then all non-colliders are blockable. So, the partial condition for  $\sigma$ -separation “a blockable non-collider in  $C$ ” can be simplified to “(any) non-collider in  $C$ ”.*

*So in the acyclic case we can simplify the notion of  $\sigma$ -separation, which is usually referred to as  $d$ -separation. However, in the non-acyclic setting  $d$ -separation (“(any) non-collider in  $C$ ”) and  $\sigma$ -separation (“a blockable non-collider in  $C$ ”) are clearly not equivalent anymore.*

*It turned out that in the non-acyclic case  $\sigma$ -separation is the more general concept (and as said above it also captures the acyclic case equivalently well), see [FM17, FM18, FM20, BFPM21]. We will first focus on CADMGs (acyclic) for which we can restrict ourselves to the somewhat simpler  $d$ -separation. Later, we will pick up  $\sigma$ -separation again when we deal with cycles.*

### Proofs - $\sigma$ -Open Walks and Paths

The following lemma will be convenient to relate  $\sigma$ -open walks and paths in the notion of  $\sigma$ -separation.

**Lemma 3.3.10.** *Let  $G = (J, V, E, L)$  be a CDMG,  $C \subseteq V \cup J$  and  $\pi = (v_0 \ast\ast \dots \ast\ast v_n)$  be a  $C$ - $\sigma$ -open walk in  $G$ . Suppose  $v_i \in \text{Sc}^G(v_j)$  for some  $i, j \in \{0, \dots, n\}$  with  $i < j$ . If we then replace the subwalk  $v_i \ast\ast \dots \ast\ast v_j$  of  $\pi$  by*

- (i) *a shortest directed path  $v_i \rightarrow \dots \rightarrow v_j$  in  $G$  if  $j = n$  or if  $v_j \rightarrow v_{j+1}$  on  $\pi$ , or*
- (ii) *a shortest directed path  $v_i \leftarrow \dots \leftarrow v_j$  in  $G$  otherwise,*

*then this new subwalk is entirely within  $\text{Sc}^G(v_j)$  and the modified walk  $\pi'$  is still  $C$ - $\sigma$ -open.*

*Proof.*  $\pi'$  cannot become  $C$ - $\sigma$ -blocked at one of the initial nodes  $v_0, \dots, v_{i-1}$  or at one of the final nodes  $v_{j+1}, \dots, v_n$  on  $\pi'$ , since these nodes occur in the same local configuration on  $\pi$  and are not  $C$ - $\sigma$ -blocked on  $\pi$  by assumption. Furthermore,  $\pi'$  cannot become  $C$ - $\sigma$ -blocked at one of the nodes strictly between  $v_i$  and  $v_j$  on  $\pi'$  (if there are any), since these nodes are all non-endnode non-colliders that only point to nodes in the same strongly connected component  $\text{Sc}^G(v_j)$ . It is also worth noting that  $\pi'$  cannot become  $C$ - $\sigma$ -blocked at any of its endnodes, which could be  $v_i$  or  $v_j$  or both, because those are the same in  $\pi$ . So in the following we can w.l.o.g. assume that both  $v_i$  and  $v_j$  are non-endnodes of  $\pi$  and thus  $\pi'$ .

Case (i). By assumption  $v_j$  is either a fork or a right chain (or the right endnode) on  $\pi$  that is  $C$ - $\sigma$ -open. Since the same blocking criteria apply to  $v_j$  on  $\pi'$  it remains  $C$ - $\sigma$ -open on  $\pi'$ . If  $v_i = v_j$  then also  $v_i$  is  $C$ - $\sigma$ -open on  $\pi'$  (if  $v_i$  is the left endnode or not). If  $v_i \neq v_j$ , then the new directed path  $v_i \rightarrow \dots \rightarrow v_j$  in  $\pi'$  is  $C$ - $\sigma$ -open at  $v_i$  because all nodes in between lie in the same strongly connected component  $\text{Sc}^G(v_i)$  (or  $v_i$  is the left endnode anyways).

Case (ii). Since case (i) is solved we can assume that we have  $j < n$  with  $v_j \leftarrow\ast v_{j+1}$  in  $\pi$ . If  $v_{i-1} \leftarrow\ast v_i$  on  $\pi'$  (or  $v_i$  the left endnode) then this case is analogous to case (i). So we can also assume that we have  $i > 0$  and  $v_{i-1} \ast\rightarrow v_i$  on  $\pi$ . So  $\pi$  looks as follows:

$$\pi : \quad \dots v_{i-1} \ast\rightarrow v_i \ast\ast \dots \ast\ast v_j \leftarrow\ast v_{j+1} \dots$$



So there must be a smallest number  $k \in \{i, \dots, j\}$  such that a collider appears at  $v_k$  on  $\pi$ :

$$\pi : \quad \dots v_{i-1} \ast \rightarrow v_i \rightarrow \dots \rightarrow v_k \leftarrow \ast \dots \ast \rightarrow v_j \leftarrow \ast v_{j+1} \dots .$$

Since  $\pi$  is  $C$ - $\sigma$ -open we have  $v_k \in \text{Anc}^G(C)$ . Since  $v_i \in \text{Anc}^G(v_k)$  (otherwise  $v_k$  would not be the first collider appearing after  $v_i$ ) we thus have that also  $v_i \in \text{Anc}^G(C)$ . So if we replace the subwalk  $v_i \ast \rightarrow \dots \ast \rightarrow v_j$  of  $\pi$  by the shortest directed path  $v_i \leftarrow \dots \leftarrow v_j$  in  $G$  we then get for  $\pi'$  the following situation:

$$\pi' : \quad \dots v_{i-1} \ast \rightarrow v_i \leftarrow \dots \leftarrow v_j \leftarrow \ast v_{j+1} \dots ,$$

which is then  $C$ - $\sigma$ -open at  $v_i$  as  $v_i \in \text{Anc}^G(C)$ . Note that this holds also when  $v_i = v_j$ . If  $v_i \neq v_j$  then  $v_j$  is also  $C$ - $\sigma$ -open on  $\pi'$  as  $v_j$  points left to a node in the same strongly connected component as  $v_j$ .

So in all cases  $\pi'$  stays  $C$ - $\sigma$ -open.  $\square$

**Proposition 3.3.6.** *Let  $G = (J, V, E, L)$  be a CDMG. For  $C \subseteq J \cup V$ , and  $w_1, w_2 \in J \cup V$ , the following are equivalent:*

1. *there exists a  $C$ - $\sigma$ -open path between  $w_1$  and  $w_2$  in  $G$ ;*
2. *there exists a  $C$ - $\sigma$ -open walk between  $w_1$  and  $w_2$  in  $G$ ;*
3. *there exists a  $C$ - $\sigma$ -open walk between  $w_1$  and  $w_2$  in  $G$  such that all its colliders lie in  $C$  (and not just in  $\text{Anc}^G(C)$ ).*

*Proof.* 3  $\implies$  2 and 1  $\implies$  2 are trivial. Note that paths are walks.

2  $\implies$  3: Suppose there exists a  $C$ - $\sigma$ -open walk  $\pi$  from  $w_1$  to  $w_2$ . Then consider a collider  $v_{k-1} \ast \rightarrow v_k \leftarrow \ast v_{k+1}$  on  $\pi$  with  $v_k \in \text{Anc}^G(C) \setminus C$ . So there exists a non-trivial directed path from  $v_k$  to a node  $c_k \in C$  with all other nodes not in  $C$ . If we then replace the collider at  $v_k$  in  $\pi$  by that path and its reverse we get:

$$\dots \ast \rightarrow v_{k-1} \ast \rightarrow v_k \rightarrow \dots \rightarrow c_k \leftarrow \dots \leftarrow v_k \leftarrow \ast v_{k+1} \ast \rightarrow \dots .$$

This walk is then  $C$ - $\sigma$ -open at all places between  $v_k$  on the left and  $v_k$  on the right because they are non-colliders not in  $C$ . If we do this iteratively for all colliders not in  $C$  we get the desired  $C$ - $\sigma$ -open walk where all colliders lie in  $C$ .

2  $\implies$  1: Let  $\pi = (v_0 \ast \rightarrow \dots \ast \rightarrow v_n)$  be a  $C$ - $\sigma$ -open walk between nodes  $v_0 = w_1$  and  $v_n = w_2$  in  $G$ . If a node  $w$  occurs more than once on  $\pi$ , let  $v_i$  be the first node on  $\pi$  and  $v_j$  be the last node on  $\pi$  that are in  $\text{Sc}^G(w)$ . We now use Lemma 3.3.10 to construct a new walk  $\pi'$  from  $\pi$  by replacing the subwalk between  $v_i$  and  $v_j$  of  $\pi$  by a particular directed path in  $\text{Sc}^G(w)$  between  $v_i$  and  $v_j$  in such a way that  $\pi'$  is still  $C$ - $\sigma$ -open. In  $\pi'$ , the number of nodes that occurs more than once is at least one less than in  $\pi$ , and all nodes within  $\text{Sc}^G(w)$  occur within a single segment. This replacement procedure can be repeated until no nodes occur more than once. We have then obtained a  $C$ - $\sigma$ -open path between  $w_1$  and  $w_2$ .  $\square$

### 3.4. d-Separation

**Definition 3.4.1** (d-blocked walks). Let  $G = (J, V, E, L)$  be a CDMG and  $C \subseteq J \cup V$  a subset of nodes and  $\pi$  a walk in  $G$ :

$$\pi = (v_0 \ast\ast \cdots \ast\ast v_n).$$

1. We say that the walk  $\pi$  is  $C$ -d-blocked or d-blocked by  $C$  to emphasize the use of the bi-directed edges.<sup>14</sup> if either:

i.)  $v_0 \in C$  or  $v_n \in C$  or:

ii.) there are two adjacent edges in  $\pi$  of one of the following forms:

$$\begin{array}{llll} \text{left chain:} & v_{k-1} \leftarrow v_k \leftarrow\ast v_{k+1} & \text{with} & v_k \in C, \\ \text{right chain:} & v_{k-1} \ast\rightarrow v_k \rightarrow v_{k+1} & \text{with} & v_k \in C, \\ \text{fork:} & v_{k-1} \leftarrow v_k \rightarrow v_{k+1} & \text{with} & v_k \in C, \\ \text{collider:} & v_{k-1} \ast\rightarrow v_k \leftarrow\ast v_{k+1} & \text{with} & v_k \notin \text{Anc}^G(C). \end{array}$$

2. We say that the walk  $\pi$  is  $C$ -d-open if it is not  $C$ -d-blocked.

**Remark 3.4.2.** If we consider end-nodes, left chains, right chains and forks as non-colliders then we can simply state:

$\pi$  is d-blocked by  $C$  if and only if it either contains a non-collider in  $C$  or a collider not in  $\text{Anc}^G(C)$ .

**Definition 3.4.3** (d-separation). Let  $G = (J, V, E, L)$  be a CDMG and  $A, B, C \subseteq J \cup V$  (not necessarily disjoint) subset of nodes. We then say that:

1.  $A$  is d-separated from  $B$  given  $C$  in  $G$ , in symbols:

$$A \perp_G^d B \mid C,$$

if every walk from a node in  $A$  to a node in  $J \cup B$  (sic!)<sup>13</sup> is  $C$ -d-blocked by  $C$ .

2. If that property does not hold we will write:

$$A \not\perp_G^d B \mid C.$$

**Remark 3.4.4.** 1. A similar result from Proposition 3.3.6 holds for d-separation as well.

2. d-separation is stable under marginalization, similar to Lemma 3.3.8.

<sup>14</sup>The “d” here stands for “directional”. d-separation was first only used for DAGs (without bi-directed edges). For ADMGs it was then called m-separation in [Ric03] But since the notion of m-separation is arguably the natural extension of d-separation and to avoid introducing more definitions, we will just call it d-separation as well, which will not create any ambiguity.

### 3.5. Acyclifications

It is possible to reformulate the notion of  $\sigma$ -separation in terms of  $d$ -separation on a modified and acyclic graph by making use of the following construction, which will be the main tool to extend the acyclic theory to the cyclic one. The construction was first proposed in the context of CBNs by [Spi94, Spi95].

**Definition 3.5.1.** *Given a CDMG  $G = (J, V, E, L)$ , we call a CADMG  $G' = (J', V', E', L')$  an acyclification of  $G$  if*

- (i)  $G'$  is acyclic;
- (ii)  $G'$  has the same input nodes and output nodes as  $G$ , i.e.  $J' = J$  and  $V' = V$ ;
- (iii) for every pair of nodes  $(i, j)$  such that  $i \notin \text{Sc}^G(j)$ :
  - a)  $i \rightarrow j \in E'$  iff there exists a node  $j' \in \text{Sc}^G(j)$  such that  $i \rightarrow j' \in E$ ;
  - b)  $i \leftrightarrow j \in L'$  iff there exist nodes  $i' \in \text{Sc}^G(i), j' \in \text{Sc}^G(j)$  such that  $i' \leftrightarrow j' \in L$ ;
- (iv) for every pair of distinct nodes  $(i, j)$  such that  $i \in \text{Sc}^G(j)$ :  $i \rightarrow j \in E'$  or  $i \leftarrow j \in E'$  or  $i \leftrightarrow j \in L'$ .

The important property of acyclifications is that they can be used to express  $\sigma$ -separation in a (possibly cyclic) graph in terms of  $d$ -separation in an acyclification.

**Proposition 3.5.2.** *Let  $G = (J, V, E, L)$  be a CDMG and  $G'$  an acyclification of  $G$ . Then for  $A, B, C \subseteq V \cup J$  (not necessarily disjoint) subsets of nodes we have the equivalence:*

$$A \perp_G^\sigma B | C \iff A \perp_{G'}^\sigma B | C \iff A \perp_{G'}^d B | C.$$

*Proof.* We will show that there is a  $C$ - $\sigma$ -open walk between  $A$  and  $B \cup J$  in  $G$  if and only if there is a  $C$ - $\sigma$ -open walk between  $A$  and  $B \cup J$  in  $G'$ . Since  $G'$  is acyclic, this is in turn equivalent to the existence of a  $C$ - $d$ -open walk between  $A$  and  $B \cup J$  in  $G'$ .

$\implies$  : Suppose there is a  $C$ - $\sigma$ -open walk  $\pi = (v_0, \dots, v_n)$  between  $A$  and  $B \cup J$  in  $G$ . All its colliders are in  $C$  and all its non-colliders are either not in  $C$ , or otherwise, point only to nodes in the same strongly connected component. Note that each edge between two nodes in different strongly connected components in  $G$  is also present in  $G'$ . Edges between two nodes in the same strongly connected component, however, may not be present in  $G'$ . Therefore, we will replace these edges with walks in  $G'$ . Consider a subwalk  $(v_i, \dots, v_j)$  of maximum length that is entirely contained within a strongly connected component in  $G$  (with possibly  $i = j$ ). We distinguish different cases and show for each case how this subwalk can be replaced by a subwalk in  $G'$ .

- (i)  $\ast \rightarrow v_i \dots v_j \leftarrow \ast$ : the subwalk between  $v_i$  and  $v_j$  has to contain a collider, say  $w$ , which must be in  $C$  since the walk between  $v_i$  and  $v_j$  is  $C$ - $\sigma$ -open. We can replace this subwalk by  $\ast \rightarrow w \leftarrow \ast$  in  $G'$  such that  $w$  becomes a collider in  $C$ .

- (ii)  $(\leftarrow)v_i \cdots v_j \leftarrow*$ :<sup>15</sup> here  $v_i$  is a non-collider pointing to another strongly connected component or  $v_i$  is an endnode, and in both cases,  $v_i \notin C$ . Therefore, we can replace the subwalk by  $(\leftarrow)v_i \leftarrow*$  in  $G'$ , such that  $v_i$  becomes a non-collider not in  $C$ .
- (iii)  $*\rightarrow v_i \cdots v_j(\rightarrow)$ : analogous to the previous case, we can replace it by  $*\rightarrow v_j(\rightarrow)$  in  $G'$ , such that  $v_j$  becomes a non-collider not in  $C$ .
- (iv)  $(\leftarrow)v_i \cdots v_j(\rightarrow)$ :  $v_i, v_j$  are both not in  $C$  by assumption. If  $i = j$ , we replace this subwalk by  $(\leftarrow)v_i(\rightarrow)$  such that  $v_i$  becomes a non-collider not in  $C$ . If  $i < j$ , we replace this subwalk by  $(\leftarrow)v_i ** v_j(\rightarrow)$  with  $v_i ** v_j$  any edge connecting  $v_i$  and  $v_j$  in  $G'$ , such that both  $v_i$  and  $v_j$  become non-colliders not in  $C$ .

By replacing all maximal subwalks of the original walk  $\pi$  that are contained within a single strongly connected component of  $G$  in this way, we obtain a walk in the acyclification  $G'$  that is  $C$ - $\sigma$ -open by construction. Note that the modified walk has the same endpoints ( $v_0$  and  $v_n$ ) as the original walk.

$\Leftarrow$  : Suppose there is a  $C$ - $\sigma$ -open walk  $\pi'$  between  $A$  and  $B \cup J$  in  $G'$ . All its colliders are in  $C$ , and all its non-colliders are not in  $C$ . We will construct a walk  $\pi$  in  $G$  with the same endpoints as  $\pi'$  that is  $C$ - $\sigma$ -open.

Consider a non-trivial subwalk  $(v_i, \dots, v_j)$  on  $\pi'$  of maximum length that is entirely contained within a strongly connected component of  $G$ . This subwalk may not be present in  $G'$ . We distinguish different cases and show for each case how this subwalk can be replaced by a subwalk in  $G$ .

- (i)  $*\rightarrow v_i \cdots v_j \leftarrow*$ : the subwalk between  $v_i$  and  $v_j$  has to contain a collider, say  $w$ , which must be in  $C$  since the walk between  $v_i$  and  $v_j$  is  $C$ - $\sigma$ -open, and must be in  $\text{Sc}^G(v_i) = \text{Sc}^G(v_j)$  by assumption. We can replace this subwalk by  $*\rightarrow v_i \rightarrow \dots \rightarrow w \leftarrow \dots \leftarrow v_j \leftarrow*$  in  $G$ , with possibly  $v_i = w$  and possibly  $w = v_j$ , with all nodes in  $\text{Sc}^G(v_i)$ . Note that the modified walk remains  $C$ - $\sigma$ -open.
- (ii)  $(\leftarrow)v_i \cdots v_j \leftarrow*$ : here  $v_i$  is a non-collider pointing to another strongly connected component or  $v_i$  is an endnode, and in both cases,  $v_i \notin C$ . We can replace this subwalk by a shortest directed walk  $(\leftarrow)v_i \leftarrow \dots \leftarrow v_j \leftarrow*$  in  $G$  with all nodes in  $\text{Sc}^G(v_i)$ . Note that the modified walk remains  $C$ - $\sigma$ -open.
- (iii)  $*\rightarrow v_i \cdots v_j(\rightarrow)$ : analogous to the previous case, we can replace it by  $*\rightarrow v_i \rightarrow \dots \rightarrow v_j(\rightarrow)$  in  $G$ .
- (iv)  $(\leftarrow)v_i \cdots v_j(\rightarrow)$ :  $v_i, v_j$  are both not in  $C$  by assumption. We can replace this subwalk by a shortest directed walk  $(\leftarrow)v_i \rightarrow \dots \rightarrow v_j(\rightarrow)$  in  $G$  with all nodes in  $\text{Sc}^G(v_i)$ . The modified walk remains  $C$ - $\sigma$ -open.

---

<sup>15</sup>We put parentheses around the first directed edge to indicate that this case also applies if  $v_i$  is an endnode, i.e., if  $i = 0$ .

In each of the four cases, in the modified walk both  $v_i$  and  $v_j$  become either colliders in  $C$ , or non-colliders not in  $C$ , or non-colliders in  $C$  that only point to a node in the same strongly connected component of  $G$ .

Now, we will replace edges on  $\pi'$  between two strongly connected components that are not present in  $G$ . For any directed edge  $i \rightarrow j$  on  $\pi'$  with  $j \notin \text{Sc}^G(i)$ , there must be a  $j' \in \text{Sc}^G(j)$  such that  $i \rightarrow j'$  is present in  $G$ , and hence there must be a directed path  $j' \rightarrow \dots \rightarrow j$  entirely in  $\text{Sc}^G(j)$  such that we can use  $i \rightarrow j' \rightarrow \dots \rightarrow j$  as replacement in  $G$  of the edge  $i \rightarrow j$ . Similarly, for any bidirected edge  $i \leftrightarrow j$  on  $\pi'$  with  $j \notin \text{Sc}^G(i)$ , there must be  $i' \in \text{Sc}^G(i)$  and  $j' \in \text{Sc}^G(j)$  such that  $i' \leftrightarrow j'$  is present in  $G$ , and hence there must be a walk  $i \leftarrow \dots \leftarrow i' \leftrightarrow j' \rightarrow \dots \rightarrow j$  in  $G$ , where  $i \leftarrow \dots \leftarrow i'$  is entirely in  $\text{Sc}^G(i)$  and  $j' \rightarrow \dots \rightarrow j$  is entirely in  $\text{Sc}^G(j)$ , that we can use as replacement in  $G$  of the edge  $i \leftrightarrow j$ . The new nodes introduced on  $\pi$  in these replacements are non-colliders that only point to nodes in the same strongly connected component. The endpoints of the replacement paths do not change their status: if they were colliders in  $C$  on  $\pi'$  they still are on  $\pi$ , and if they were non-colliders not in  $C$  on  $\pi'$  they still are on  $\pi$ .

Hence we have constructed a walk  $\pi$  in  $G$  with the same endpoints as  $\pi'$  that is  $C$ - $\sigma$ -open.  $\square$

The following construction shows that acyclifications exist (but it is just one out of many possible ways to construct acyclifications).

**Example 3.5.3** (The standard acyclification). *Let  $G = (J, V, E, L)$  be a CDMG. Then we define the standard acyclification of  $G$  as the CDMG  $G' = (J, V, E', L')$  where:*

$$\begin{aligned} E' &:= \{v_1 \rightarrow v_2 \mid v_1 \in J \cup V, v_2 \in V, v_2 \notin \text{Sc}^G(v_1), \exists v'_2 \in \text{Sc}^G(v_2) : v_1 \rightarrow v'_2 \in E\}, \\ L' &:= \{v_1 \leftrightarrow v_2 \mid v_1, v_2 \in V, v_1 \neq v_2, \exists v'_1 \in \text{Sc}^G(v_1), v'_2 \in \text{Sc}^G(v_2) : v_1 \leftrightarrow v'_2 \in L\} \\ &\cup \{v_1 \leftrightarrow v_2 \mid v_1, v_2 \in V, v_1 \neq v_2, v_1 \in \text{Sc}^G(v_2)\}. \end{aligned}$$

*The standard acyclification of a CDMG is acyclic, i.e. a CADMG.*

*Proof.* Assume that  $G'$  is not acyclic. Then there exists a non-trivial cyclic directed walk in  $G'$ :

$$v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_k \rightarrow v_1,$$

for some  $k \geq 1$ . It is clear that  $k \geq 2$  because clearly  $v_1 \in \text{Sc}^G(v_1)$ , which rules out the existence of an edge  $v_1 \rightarrow v_1 \in G'$ . For simplicity we now identify  $v_{k+1} := v_1$  in the following. By construction of  $G'$  for every  $i = 1, \dots, k$  there exists  $v'_{i+1} \in \text{Sc}^G(v_{i+1})$  such that the edge  $v_i \rightarrow v'_{i+1}$  exists in  $G$ . Since  $v'_{i+1} \in \text{Sc}^G(v_{i+1})$  there also exists a directed walk in  $G$ :

$$v'_{i+1} \rightarrow \dots \rightarrow v_{i+1}.$$

Concatenating all directed walks we get the cyclic directed walk in  $G$ :

$$v_1 \rightarrow v'_2 \rightarrow \dots \rightarrow v_2 \rightarrow v'_3 \rightarrow \dots \rightarrow v_k \rightarrow v'_1 \rightarrow \dots \rightarrow v_1.$$

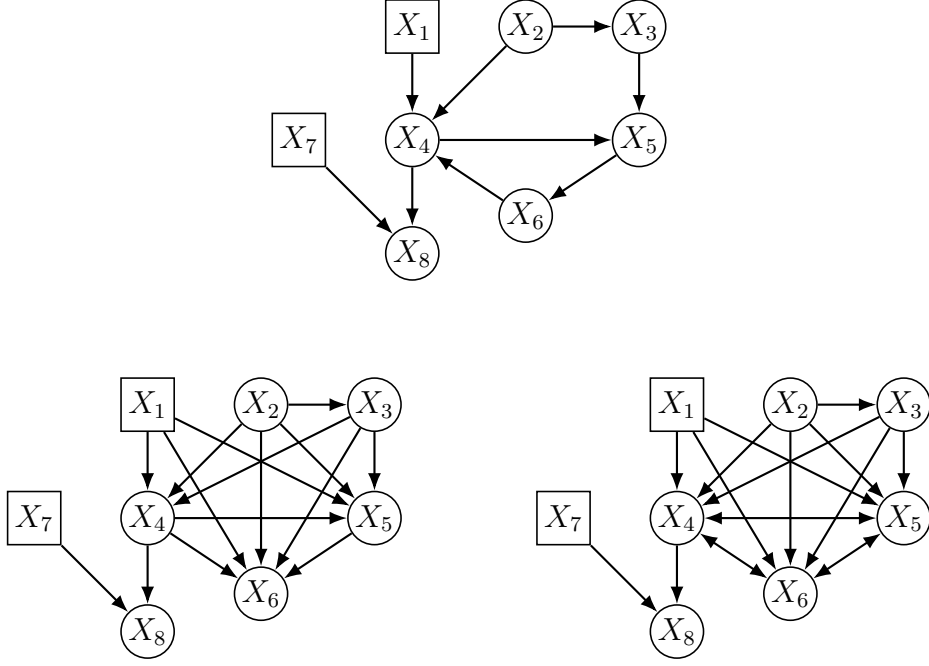


Figure 7: Top: CDMG  $G$ . Bottom: two acyclifications of  $G$ .

This shows that  $v_2 \in \text{Sc}^G(v_1)$ , which is a contradiction to the existence of the edge  $v_1 \rightarrow v_2 \in G'$ . So a non-trivial cyclic directed walk in  $G'$  cannot exist in the first place. So  $G'$  must be acyclic.  $\square$

### 3.6. Separoid Axioms for $\sigma$ -/d-Separation

**Definition/Theorem 3.6.1** ((Asymmetric) separoid axioms for  $\sigma$ -separation/d-separation).

Let  $G = (J, V, E, L)$  be a CDMG and  $A, B, C, D \subseteq J \cup V$  subsets of nodes. Then the ternary relations  $\perp = \perp_G^d$  and  $\perp = \perp_G^\sigma$  satisfy the following rules:

a) *Extended Left Redundancy:*

$$D \subseteq A \implies D \perp B | A.$$

b) *J-Restricted Right Redundancy:*

$$A \perp \emptyset | C \cup J \text{ always holds.}$$

c) *J-Inverted Right Decomposition:*

$$A \perp B | C \implies A \perp J \cup B | C.$$

d) *Left Decomposition:*

$$A \cup D \perp B | C \implies D \perp B | C.$$

e) *Right Decomposition:*

$$A \perp B \cup D | C \implies A \perp D | C.$$

f) *Left Weak Union:*

$$A \cup D \perp B | C \implies A \perp B | D \cup C.$$

g) *Right Weak Union:*

$$A \perp B \cup D | C \implies A \perp B | D \cup C.$$

h) *Left Contraction:*

$$(A \perp B | D \cup C) \wedge (D \perp B | C) \implies A \cup D \perp B | C.$$

i) *Right Contraction:*

$$(A \perp B | D \cup C) \wedge (A \perp D | C) \implies A \perp B \cup D | C.$$

j) *Right Cross Contraction:*

$$(A \perp B | D \cup C) \wedge (D \perp A | C) \implies A \perp B \cup D | C.$$

k) *Flipped Left Cross Contraction:*

$$(A \perp B | D \cup C) \wedge (B \perp D | C) \implies B \perp A \cup D | C.$$

*In particular, we have the equivalences:*

$$(A \perp B \cup D | C) \iff (A \perp B | D \cup C) \wedge (A \perp D | C),$$

$$(A \cup D \perp B | C) \iff (A \perp B | D \cup C) \wedge (D \perp B | C).$$

*We also get:*

l) *J-Restricted Symmetry:*

$$A \perp B | C \cup J \implies B \perp A | C \cup J.$$

*For the special case of  $J = \emptyset$  we have thus (unrestricted) Symmetry.*

**Remark 3.6.2.** *Let the assumptions be like in Theorem 3.6.1. We also have the following rules:*

m) *Left Composition:*

$$(A \perp B | C) \wedge (D \perp B | C) \implies A \cup D \perp B | C.$$

n) *Right Composition:*

$$(A \perp B | C) \wedge (A \perp D | C) \implies A \perp B \cup D | C.$$

o) *Left Intersection: If  $A \cap D = \emptyset$  then:*

$$(A \perp B \mid D \cup C) \wedge (D \perp B \mid A \cup C) \implies A \cup D \perp B \mid C.$$

p) *Right Intersection: If  $B \cap D = \emptyset$  then:*

$$(A \perp B \mid D \cup C) \wedge (A \perp D \mid B \cup C) \implies A \perp B \cup D \mid C.$$

### Proofs - Separoid Axioms for $\sigma$ -/d-Separation

In the following let  $G = (J, V, E, L)$  be a CDMG and  $A, B, C, D \subseteq J \cup V$  (not necessarily disjoint) subsets of nodes.

Recall that we say that  $A$  is  $\sigma$ -separated from  $B$  given  $C$  in  $G$ , in symbols:

$$A \underset{G}{\perp}^{\sigma} B \mid C,$$

if every walk from a node in  $A$  to a node in  $J \cup B$  (sic!) is  $\sigma$ -blocked by  $C$ .

Again, a walk  $\pi$  is  $\sigma$ -blocked by  $C$  if it either contains a blockable non-collider in  $C$  or a collider not in  $C$ .

We abbreviate the ternary relations in the following as:  $\perp := \underset{G}{\perp}^{\sigma}$ .

**Lemma 3.6.3** (Extended Left Redundancy).

$$D \subseteq A \implies D \perp B \mid A.$$

*Proof.* If  $\pi$  is a walk from a node  $v$  in  $D$  to a node  $w$  in  $J \cup B$  then its first end node is in  $A$ , so  $\pi$  is  $\sigma$ -blocked by  $A$ .  $\square$

**Lemma 3.6.4** ( $J$ -Restricted Right Redundancy).

$$A \perp \emptyset \mid C \cup J \quad \text{always holds.}$$

*Proof.* If  $\pi$  is a walk from a node  $v$  in  $A$  to a node  $w$  in  $J$  then its last end node is in  $C \cup J$ , so  $\pi$  is  $\sigma$ -blocked by  $C \cup J$ .  $\square$

**Lemma 3.6.5** ( $J$ -Inverted Right Decomposition).

$$A \perp B \mid C \implies A \perp J \cup B \mid C.$$

*Proof.* If  $\pi$  is a walk from a node  $v$  in  $A$  to a node  $w$  in  $J \cup J \cup B$  then  $w \in J \cup B$ . If  $w \in J \cup B$  then by assumption  $\pi$  is  $\sigma$ -blocked by  $C$ .  $\square$

**Lemma 3.6.6** (Left Decomposition).

$$A \cup D \perp B \mid C \implies D \perp B \mid C.$$

*Proof.* If  $\pi$  is a walk from a node  $v$  in  $D$  to a node  $w$  in  $J \cup B$ , then the walk  $\pi$  is also a walk from  $A \cup D$  to  $J \cup B$ , which by assumption is  $\sigma$ -blocked by  $C$ .  $\square$



**Lemma 3.6.7** (Right Decomposition).

$$A \perp B \cup D | C \implies A \perp D | C.$$

*Proof.* If  $\pi$  is a walk from a node  $v$  in  $A$  to a node  $w$  in  $J \cup D$ , then the walk  $\pi$  is also a walk from  $A$  to  $J \cup B \cup D$ , which by assumption is  $\sigma$ -blocked by  $C$ .  $\square$

**Lemma 3.6.8** (Left Weak Union).

$$A \cup D \perp B | C \implies A \perp B | D \cup C.$$

*Proof.* Let us assume the contrary:  $A \not\perp B | D \cup C$ . Then there exists a shortest  $(D \cup C)$ - $\sigma$ -open walk  $\pi$  from a node  $v$  in  $A$  to a node  $w$  in  $J \cup B$  in  $G$  such that every collider of  $\pi$  is in  $D \cup C$ . Then every blockable non-collider of  $\pi$  is not in  $D \cup C$ .

If now  $\pi$  does not contain any node from  $D \setminus C$  then every collider of  $\pi$  is in  $C$ . This implies that  $\pi$  is  $C$ - $\sigma$ -open, which contradicts the assumption:  $A \cup D \perp B | C$ .

So we can assume now that  $\pi$  contains a node in  $D \setminus C$ . Then consider the shortest sub-walk  $\tilde{\pi}$  in  $\pi$  starting from the end-node  $w \in J \cup B$  and going back to the first node  $u \in D \setminus C$ . This means that  $\tilde{\pi}$  is a walk from  $D \setminus C$  to  $J \cup B$  where the end-node  $u$  of  $\tilde{\pi}$  is the only node in  $D \setminus C$ . So  $\tilde{\pi}$  does not contain any collider in  $D \setminus C$ . So all colliders of  $\tilde{\pi}$  lie in  $C$ . All blockable non-colliders of  $\tilde{\pi}$  that are different from the end-node  $u$  are also blockable non-colliders on  $\pi$ . They are thus not in  $D \cup C$  by the assumption on  $\pi$ , in particular, not in  $C$ . The only remaining blockable non-collider  $u$  of  $\tilde{\pi}$  lies in  $D \setminus C$  by construction and it thus lies not in  $C$  either. So  $\tilde{\pi}$  is  $C$ - $\sigma$ -open walk from  $A \cup D$  to  $J \cup B$ . This contradicts the assumption:  $A \cup D \perp B | C$ .

So the premise:  $A \not\perp B | D \cup C$ , must be false. This shows:  $A \perp B | D \cup C$ .  $\square$

**Lemma 3.6.9** (Right Weak Union).

$$A \perp B \cup D | C \implies A \perp B | D \cup C.$$

*Proof.* Follow the same steps as in Left Weak Union (Lemma 3.6.8). Soe there exists a shortest  $(D \cup C)$ - $\sigma$ -open walk  $\pi$  from a node  $v$  in  $A$  to a node  $w$  in  $J \cup B \cup D$  in  $G$  such that every collider of  $\pi$  is in  $D \cup C$ . If  $\pi$  does not contain any nodes from  $D \setminus C$  we get a contradiction to:  $A \perp B \cup D | C$ . Then, again, we can assume that  $\pi$  contains a node in  $D \setminus C$ . Then consider the shortest sub-walk  $\tilde{\pi}$  in  $\pi$  from  $v \in A$  to a node  $u \in D \setminus C$ . This means that  $\tilde{\pi}$  does not contain any collider in  $D \setminus C$ , so they are all in  $C$ . Furthermore, all blockable non-colliders are not in  $C$ . So  $\tilde{\pi}$  is  $C$ - $\sigma$ -open walk from  $A$  to  $J \cup B \cup D$ . This contradicts the assumption:  $A \perp B \cup D | C$ .  $\square$

**Lemma 3.6.10** (Left Contraction).

$$(A \perp B | D \cup C) \wedge (D \perp B | C) \implies A \cup D \perp B | C.$$

*Proof.* Let us assume the contrary:  $A \cup D \not\perp B | C$ . Then there exists a shortest  $C$ - $\sigma$ -open walk  $\pi$  from a node  $v$  in  $A \cup D$  to a node  $w$  in  $J \cup B$  in  $G$  such that every collider of  $\pi$  is in  $C$ . So every blockable non-collider is not in  $C$ . In particular,  $v \notin C$ . We now

claim that  $v$  is the only node of  $\pi$  that is in  $(A \cup D) \setminus C$ . Otherwise, there would be a non-end-node  $u$  of  $\pi$  with  $u \in (A \cup D) \setminus C$ . Since  $u \notin C$  the whole sub-walk from  $u$  to  $w$  would already be a  $C$ - $\sigma$ -open walk from  $A \cup D$  to  $J \cup B$ , which is also shorter than  $\pi$ , which contradicts the assumption. So we can assume that  $v$  is the only node of  $\pi$  that is in  $(A \cup D) \setminus C$ . In particular, all blockable non-colliders of  $\pi$  that are different from  $v$  are not in  $D \setminus C$  and thus are not in  $D \cup C = (D \setminus C) \cup C$ .

Furthermore,  $v$  cannot lie in  $D \setminus C$  as it would contradict the assumption:  $D \perp B | C$ . It follows that  $v \in A \setminus C$  and  $\pi$  is a walk from  $A$  to  $J \cup B$  whose colliders are in  $C \subseteq D \cup C$  and all blockable non-colliders are not in  $D \cup C$ . But this contradicts the other assumption:  $A \perp B | D \cup C$ .  $\square$

**Lemma 3.6.11** (Right Contraction).

$$(A \perp B | D \cup C) \wedge (A \perp D | C) \implies A \perp B \cup D | C.$$

*Proof.* Let us assume the contrary:  $A \not\perp B \cup D | C$ . Then there exists a shortest  $C$ - $\sigma$ -open walk  $\pi$  from a node  $v$  in  $A$  to a node  $w$  in  $J \cup B \cup D$  in  $G$  such that every collider of  $\pi$  is in  $C$ . So every blockable non-collider is not in  $C$  and  $w$  is the only node of  $\pi$  that is in  $(J \cup B \cup D) \setminus C$  (otherwise  $\pi$  could be shortened).

Also  $w$  cannot lie in  $D \setminus C$  as it would contradict the assumption:  $A \perp D | C$ . Thus  $w \in (J \cup B) \setminus C$  and  $\pi$  is a walk from  $A$  to  $J \cup B$  whose colliders all are in  $C \subseteq D \cup C$  and all blockable non-colliders are not in  $D \cup C$ . But this contradicts the other assumption:  $A \perp B | D \cup C$ .  $\square$

**Lemma 3.6.12** (Right Cross Contraction).

$$(A \perp B | D \cup C) \wedge (D \perp A | C) \implies A \perp B \cup D | C.$$

*Proof.* Verbatim the same as Right Contraction (Lemma 3.6.11), only the first contradiction is with:  $D \perp A | C$ .  $\square$

**Lemma 3.6.13** (Flipped Left Cross Contraction).

$$(A \perp B | D \cup C) \wedge (B \perp D | C) \implies B \perp A \cup D | C.$$

*Proof.* Let us assume the contrary:  $B \not\perp A \cup D | C$ . Then there exists a shortest  $C$ - $\sigma$ -open walk  $\pi$  from a node  $v$  in  $B$  to a node  $w$  in  $J \cup A \cup D$  in  $G$  such that every collider of  $\pi$  is in  $C$ . So every blockable non-collider is not in  $C$  and  $w$  is the only node of  $\pi$  that is in  $(J \cup A \cup D) \setminus C$  (otherwise  $\pi$  could be shortened).

Also  $w$  cannot lie in  $(J \cup D) \setminus C$  as it would contradict the assumption:  $B \perp D | C$ . Thus  $w \in A \setminus C$  and the walk  $\pi$  (in reverse direction) is a walk from  $A$  to  $B$  whose colliders are all in  $C \subseteq D \cup C$  and all blockable non-colliders are not in  $D \cup C$ . But this contradicts the other assumption:  $A \perp B | D \cup C$ .  $\square$

**Lemma 3.6.14** ( $J$ -Restricted Symmetry).

$$A \perp B | C \cup J \implies B \perp A | C \cup J.$$

*Proof.* This follows from Flipped Left Cross Contraction (Lemma 3.6.13) with  $D = \emptyset$  and  $C \cup J$  in place of  $C$  together with  $J$ -Restricted Right Redundancy (Lemma 3.6.4).  $\square$

**Lemma 3.6.15** (Left Composition).

$$(A \perp B | C) \wedge (D \perp B | C) \implies A \cup D \perp B | C.$$

*Proof.* Let  $\pi$  be a walk from a node  $v$  in  $A \cup D$  to a node  $w$  in  $J \cup B$ . If  $v \in A$  then  $\pi$  is  $\sigma$ -blocked by  $C$  by assumption:  $A \perp B | C$ . If  $v \in D$  then  $\pi$  is  $\sigma$ -blocked by  $C$  by assumption:  $D \perp B | C$ .  $\square$

**Lemma 3.6.16** (Right Composition).

$$(A \perp B | C) \wedge (A \perp D | C) \implies A \perp B \cup D | C.$$

*Proof.* Let  $\pi$  be a walk from a node  $v$  in  $A$  to a node  $w$  in  $J \cup B \cup D$ . If  $w \in J \cup B$  then  $\pi$  is  $\sigma$ -blocked by  $C$  by assumption:  $A \perp B | C$ . If  $w \in J \cup D$  then  $\pi$  is  $\sigma$ -blocked by  $C$  by assumption:  $A \perp D | C$ .  $\square$

**Lemma 3.6.17** (Left Intersection). *Assume that  $A \cap D = \emptyset$ , then:*

$$(A \perp B | D \cup C) \wedge (D \perp B | A \cup C) \implies A \cup D \perp B | C.$$

*Proof.* Let us assume the contrary:  $A \cup D \not\perp B | C$ . Then there exists a shortest  $C$ - $\sigma$ -open walk  $\pi$  from a node  $v$  in  $A \cup D$  to a node  $w$  in  $J \cup B$  in  $G$  such that every collider of  $\pi$  is in  $C$ . So every blockable non-collider is not in  $C$  and  $v$  is the only node of  $\pi$  that is in  $(A \cup D) \setminus C$  (otherwise  $\pi$  could be shortened).

If  $v \in A$  then by the disjointness of  $A$  and  $D$  we have that  $v \notin D$ . Then  $\pi$  is a walk from  $A$  to  $J \cup B$  whose colliders are in  $C \subseteq D \cup C$  and all blockable non-colliders are not in  $(D \setminus C) \cup C = D \cup C$ . This contradicts the assumption:  $A \perp B | D \cup C$ .

If  $v \in D$  then similarly we get a contradiction:  $D \perp B | A \cup C$ .  $\square$

**Lemma 3.6.18** (Right Intersection). *Assume that  $B \cap D = \emptyset$ , then:*

$$(A \perp B | D \cup C) \wedge (A \perp D | B \cup C) \implies A \perp B \cup D | C.$$

*Proof.* Let us assume the contrary:  $A \not\perp B \cup D | C$ . Then there exists a shortest  $C$ - $\sigma$ -open walk  $\pi$  from a node  $v$  in  $A$  to a node  $w$  in  $J \cup B \cup D$  in  $G$  such that every collider of  $\pi$  is in  $C$ . So every blockable non-collider is not in  $C$  and  $w$  is the only node of  $\pi$  that is in  $(J \cup B \cup D) \setminus C$  (otherwise  $\pi$  could be shortened).

If  $w \notin B$  then  $w \in J \cup D$ . In this case  $\pi$  is a walk from  $A$  to  $J \cup D$  where every collider is in  $C \subseteq B \cup C$  and all blockable non-colliders are not in  $(B \setminus C) \cup C = B \cup C$ . So  $\pi$  is a  $(B \cup C)$ - $\sigma$ -open walk from  $A$  to  $J \cup D$ . This contradicts the assumption:  $A \perp D | B \cup C$ .

If  $w \notin D$  then  $w \in J \cup B$ . In this case  $\pi$  is a walk from  $A$  to  $J \cup B$  where every collider is in  $C \subseteq D \cup C$  and all blockable non-colliders are not in  $D \cup C$ . So  $\pi$  is a  $(D \cup C)$ - $\sigma$ -open walk from  $A$  to  $J \cup B$ . This contradicts the assumption:  $A \perp B | D \cup C$ .

Since  $B \cap D = \emptyset$  there are no other cases ( $B^c \cup D^c = J \cup V$ ) and we are done.  $\square$

**Remark 3.6.19** (Proofs for the separoid axioms for  $d$ -separation). *The proofs for the separoid axioms for  $d$ -separation are verbatim the same as above if one exchanges the word “blockable non-collider” with just the word “non-collider” everywhere.*

## 4. Causal Bayesian Networks

### 4.1. Core Concepts

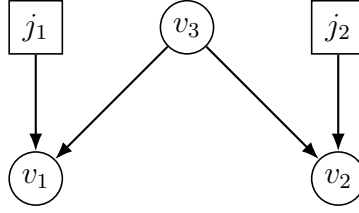


Figure 8: The Conditional Directed Acyclic Graph (CDAG) of a Causal Bayesian Network (CBN) with input variables.

**Definition 4.1.1** (Causal Bayesian network). A causal Bayesian network (CBN)—by definition—consists of:

- a conditional directed acyclic graph (CDAG):  $G = (J, V, E)$  (with finite vertex sets, and no bidirected edges),
- a standard measurable space  $\mathcal{X}_v$  for every  $v \in J \cup V$ ,
- for every  $v \in V$ , a Markov kernel:  $P_v(X_v | X_{\text{Pa}^G(v)})$ :

$$\begin{aligned} \mathcal{X}_{\text{Pa}^G(v)} &\dashrightarrow \mathcal{X}_v, \\ (A, x_{\text{Pa}^G(v)}) &\mapsto P_v(X_v \in A | X_{\text{Pa}^G(v)} = x_{\text{Pa}^G(v)}), \end{aligned}$$

where we write for  $D \subseteq J \cup V$ :

$$\begin{aligned} \mathcal{X}_D &:= \prod_{v \in D} \mathcal{X}_v, & \mathcal{X}_\emptyset &:= * = \{*\}, \\ X_D &:= (X_v)_{v \in D}, & X_\emptyset &:= *, \\ x_D &:= (x_v)_{v \in D}, & x_\emptyset &:= *. \end{aligned}$$

**Remark 4.1.2.** Most existing accounts of causal Bayesian networks do not formally distinguish input nodes from output nodes. The reasons that we do make this distinction are of a measure-theoretical nature. If all variables are discrete, and all probability mass functions and Markov kernels are strictly positive, then the formal differences between input and output nodes may be ignored and everything can be considered as output nodes.

**Definition 4.1.3** (The joint Markov kernel of a causal Bayesian network with input variables). Consider a causal Bayesian network with input variables with CDAG  $G = (J, V, E)$  with Markov kernels  $P_v(X_v | X_{\text{Pa}^G(v)})$  for  $v \in V$ . For a fixed topological ordering  $<$  of  $G$  we then define the joint Markov kernel of the CBN:

$$\mathcal{X}_J \dashrightarrow \mathcal{X}_V$$

as follows:

$$P(X_V | \text{do}(X_J)) := \bigotimes_{v \in V}^> P_v(X_v | X_{\text{Pa}^G(v)}),$$

where the nodes  $v$  run through  $V$  in reverse ordering of  $<$ , i.e. all parents are on the right of all their children.

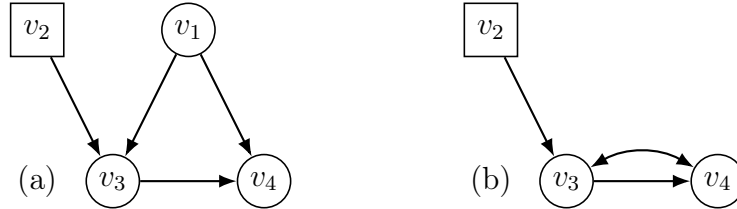


Figure 9: (a) Conditional Directed Acyclic Graph (CDAG)  $G$ ; (b) Conditional Acyclic Directed Mixed Graph (CADMG)  $G^{\setminus\{v_1\}}$  obtained after marginalizing out  $v_1$ .

**Example 4.1.4.** The CDAG  $G$  displayed in Figure 9(a) and Markov kernels  $P_1(X_1)$ ,  $P_3(X_3|X_1, X_2)$ ,  $P_4(X_4|X_1, X_3)$  give a joint Markov kernel of a CBN:

$$P(X_1, X_3, X_4 | \text{do}(X_2)) = P_4(X_4|X_1, X_3) \otimes P_3(X_3|X_1, X_2) \otimes P_1(X_1).$$

**Exercise 4.1.5.** Show that the definition of the joint Markov kernel of a CBN is independent of the topological ordering.

**Notation 4.1.6.** By abuse of notation, we will refer to the tuple:

$$M = \left( G = (J, V, E), (P_v(X_v | X_{\text{Pa}^G(v)}))_{v \in V} \right),$$

or just to the tuple:

$$M = (G, P(X_V | \text{do}(X_J)))$$

as the CBN, keeping the single Markov kernels  $P_v(X_v | X_{\text{Pa}^G(v)})$  and the spaces  $\mathcal{X}_v$  implicit.

**Remark 4.1.7** (Marginalization and conditioning). Let  $P(X_V | \text{do}(X_J))$  be the joint Markov kernel of a CBN. We can extend it to a joint Markov kernel including  $X_J$ :

$$P(X_V, X_J | \text{do}(X_J)) = P(X_V | \text{do}(X_J)) \otimes \delta(X_J | X_J).$$

For any  $A, B \subseteq J \cup V$  we then also have the marginal conditional Markov kernel:

$$P(X_A | X_B, \text{do}(X_J)),$$

which exists by theorem 2.4.16 due to the use of standard measurable spaces and is unique up to a  $P(X_B | \text{do}(X_J))$ -null set.

Furthermore, if  $C \subseteq J$  and we have:

$$X_A \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_J | X_C,$$

then we also have a Markov kernel:

$$P(X_A | \text{do}(X_C))$$

that fits into the equation:

$$P(X_A, X_C | \text{do}(X_J)) = P(X_A | \text{do}(X_C)) \otimes P(X_C | \text{do}(X_J)).$$

Note that this  $P(X_A | \text{do}(X_C))$  is unique up to a  $P(X_C | \text{do}(X_J))$ -null set. Since, further,  $P(X_C | \text{do}(X_J)) = \delta(X_C | X_J)$ , we even get that  $P(X_A | \text{do}(X_C))$  is unique (not just up to null sets). In other words, the above conditional independence states that  $P(X_A | \text{do}(X_J))$  is only dependent on the arguments from  $X_C$  and can be represented by a Markov kernel  $P(X_A | \text{do}(X_C))$ .

**Definition 4.1.8** (Causal Bayesian network with latent variables). A causal Bayesian network with latent variables (L-CBN)—*per definition*—consists of a CBN:

$$M = \left( G^+ = (J, V^+, E^+), \left( P_v(X_v | X_{\text{Pa}^{G^+}(v)}}) \right)_{v \in V^+} \right),$$

together with a disjoint decomposition of the output nodes  $V^+ = V \dot{\cup} U$  into observed nodes  $V$  and unobserved nodes  $U$ .

**Remark 4.1.9.** In the Definition 4.1.8 we make the distinction between the set of observed nodes  $V$  and that of unobserved nodes  $U$ . We could have made that distinction already earlier in the graph theory chapters and introduce CDAGs  $G^+ = (J, (V, U), E^+)$ , where we make the distinction between these node types part of the (or a new) definition. However, most of the time these sets are mathematically treated the same way and we could just consider their union  $V^+ = V \dot{\cup} U$ . Usually the distinction between  $V$  and  $U$  is only made to indicate which variables are marginalized out. Also, it often happens that one considers the same CBN in different settings, and which variables are observed and unobserved depends on the setting (for example, during training of a classifier both features and labels are observed, while during testing only features are observed). For all these reasons, we do not consider the specification of which variables are observed and which are latent part of the model.

**Notation 4.1.10.** 1. We will also often just denote a causal Bayesian network with latent variables by the tuple:

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | X_{\text{Pa}^{G^+}(v)}}) \right)_{v \in V \cup U} \right),$$

or just:

$$M = (G^+, P(X_{V \cup U} | \text{do}(X_J))).$$

2. We refer to the marginal Markov kernel of  $M$ :

$$P(X_V | \text{do}(X_J))$$

as the observable Markov kernel.

3. We call the marginalized CADMG of  $M$ :

$$G := (J, V, E, L) := (G^+) \setminus U$$

the (induced) observable CADMG.

4. We will often just refer to  $M$  as “a CBN with observed nodes  $V$ ” or “a CBN with latent nodes  $U$ ” or “a CBN with observed CADMG  $G$ ” to mean that  $M$  is a causal Bayesian network with latent variables with latent nodes  $U$  and observed nodes  $V$ .

**Example 4.1.11.** Consider again the CADG  $G$  displayed in Figure 9(a) (see also Example 4.1.4). If we assume  $v_1$  to be a latent variable (and  $v_3, v_4$  to be observed output variables), we obtain an L-CBN

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | X_{\text{Pa}^{G^+}(v)}}) \right)_{v \in V \cup U} \right),$$

with  $J = \{v_2\}$ ,  $V = \{v_3, v_4\}$ ,  $U = \{v_1\}$ . Its induced observable CADMG  $G := (G^+) \setminus \{v_1\}$  is displayed in Figure 9(b). Its observable Markov kernel is the marginal  $P(X_3, X_4 | \text{do}(X_2))$  of the Markov kernel:

$$P(X_1, X_3, X_4 | \text{do}(X_2)) = P_4(X_4 | X_1, X_2) \otimes P_2(X_3 | X_1, X_2) \otimes P_1(X_1).$$

## 4.2. Global Markov Property

**Theorem 4.2.1** (Global Markov property for causal Bayesian networks). Consider a causal Bayesian network  $M$  with observable CADMG  $G = (J, V, E, L)$  and observable Markov kernel  $P(X_V | \text{do}(X_J))$ . Then for all  $A, B, C \subseteq J \cup V$  (not necessarily disjoint) we have the implication:

$$A \perp_G^d B | C \quad \Longrightarrow \quad X_A \perp_{P(X_V | \text{do}(X_J))} X_B | X_C.$$

**Remark 4.2.2.** If one wants to make the implicit dependence on  $J$  in Theorem 4.2.1 more explicit one can equivalently also write:

$$A \perp_G^d J \cup B | C \quad \Longrightarrow \quad X_A \perp_{P(X_V | \text{do}(X_J))} X_J, X_B | X_C.$$

**Notation 4.2.3.** Let  $A, B, C \subseteq J \cup V$  with  $X_A \perp_{P(X_V | \text{do}(X_J))} X_B | X_C$ , then we have a factorization:

$$P(X_A, X_B, X_C | \text{do}(X_J)) = Q(X_A | X_C) \otimes P(X_B, X_C | \text{do}(X_J)),$$

for some Markov kernel:  $Q(X_A | X_C)$ . If we marginalize out  $X_B$  and the deterministic  $X_{C \cap J}$ , we get:

$$P(X_A, X_{C \cap V} | \text{do}(X_J)) = Q(X_A | X_C) \otimes P(X_{C \cap V} | \text{do}(X_J)).$$

So we see that  $Q(X_A|X_C)$  is a conditional Markov kernel:

$$P(X_A|X_{C \cap V}, \text{do}(X_J))$$

that does only depend on  $X_{J \cap C}$  in the do-part. So we will use the following notation for  $Q(X_A|X_C)$  (or in any other order behind the conditioning line):

$$P(X_A|X_{C \cap V}, \text{do}(X_{C \cap J}), \underline{\text{do}}(\overline{X_J})) := Q(X_A|X_C).$$

Note that by Theorem 2.5.28 we may (but do not need to) explicitly mention  $X_B$  as in:

$$P(X_A|\overline{X_B}, X_{C \cap V}, \text{do}(X_{C \cap J}), \underline{\text{do}}(\overline{X_J})),$$

because the Markov kernels are almost surely equal:

$$P(X_A|X_{C \cap V}, \text{do}(X_{C \cap J}), \underline{\text{do}}(\overline{X_J})) = P(X_A|\overline{X_B}, X_{C \cap V}, \text{do}(X_{C \cap J}), \underline{\text{do}}(\overline{X_J})) \quad P(X_C|X_J)\text{-a.s.}$$

In these suggestive notations we can state the global Markov property (Theorem 4.2.1) as:

$$\begin{aligned} A \stackrel{d}{\perp}_G B | C \\ \implies P(X_A|X_B, X_C, \text{do}(X_J)) \\ &= P(X_A|\overline{X_B}, X_{C \cap V}, \text{do}(X_{C \cap J}), \underline{\text{do}}(\overline{X_J})) \quad P(X_B, X_C|X_J)\text{-a.s.} \\ &= P(X_A|X_{C \cap V}, \text{do}(X_{C \cap J}), \underline{\text{do}}(\overline{X_J})) \quad P(X_B, X_C|X_J)\text{-a.s.} \end{aligned}$$

### Proofs - Global Markov Property

The proof of the global Markov property follows similar arguments as used in [LDLL90, Ver93, Ric03, FM17, FM18, RERS23], namely chaining the separoid axioms together in an inductive way. The main difference here is that we never rely on the Symmetry property but instead use the left and right versions of the separoid axioms separately.

**Theorem 4.2.4** (Global Markov property for causal Bayesian networks). *Consider a causal Bayesian network  $M$  with observable CADMG  $G = (J, V, E, L)$  and observable Markov kernel  $P(X_V | \text{do}(X_J))$ . Then for all  $A, B, C \subseteq J \cup V$  (not-necessarily disjoint) we have the implication:*

$$A \stackrel{d}{\perp}_G B | C \quad \implies \quad X_A \underset{P(X_V | \text{do}(X_J))}{\perp\!\!\!\perp} X_B | X_C.$$

If one wants to make the implicit dependence on  $J$  more explicit one can equivalently also write:

$$A \stackrel{d}{\perp}_G J \cup B | C \quad \implies \quad X_A \underset{P(X_V | \text{do}(X_J))}{\perp\!\!\!\perp} X_J, X_B | X_C.$$



*Proof.* Because d-separation is preserved under marginalization:

$$A \perp_G^d B | C \iff A \perp_{G^+}^d B | C,$$

we can directly assume that we work with the causal Bayesian network without latent variables that marginalizes to the given one. So w.l.o.g.  $L = \emptyset$  and  $G$  is a CDAG. We then do induction by  $\#V$ .

0.) Induction start:  $V = \emptyset$ . This means that  $A, B, C \subseteq J$ . The assumption:

$$A \perp_G^d B | C,$$

implies that we must have that  $A \subseteq C$ . Otherwise a trivial walk from  $A \subseteq J$  to  $J \cup B$  would be  $C$ -open. Since  $A, B, C \subseteq J$  we have the factorization:

$$P(X_A, X_B, X_C | \text{do}(X_J)) = \underbrace{\bigotimes_{w \in A} \delta(X_w | X_w)}_{=: Q(X_A | X_C)} \otimes \underbrace{\bigotimes_{w \in B} \delta(X_w | X_w) \otimes \bigotimes_{w \in C} \delta(X_w | X_w)}_{=: P(X_B, X_C | \text{do}(X_J))}.$$

Because  $A \subseteq C$  the Markov kernel  $Q(X_A | X_C) := \bigotimes_{w \in A} \delta(X_w | X_w)$  really is a Markov kernel from  $\mathcal{X}_C \dashrightarrow \mathcal{X}_A$ . This already shows:

$$X_A \perp_{P(X_V | \text{do}(X_J))} \parallel X_B | X_C.$$

(IND): Induction assumption: The global Markov property holds for all causal Bayesian networks (with input variables, but without latent variables and without bi-directed edges) with  $\#V < n$  (and arbitrary  $J$ ).

1.) Now assume:  $\#V = n > 0$  and  $A \perp_G^d B | C$ .

Since  $G$  is acyclic we can find a topological order  $<$  for  $G$  where the elements of  $J$  are ordered first. Let  $v \in V$  be its last element, which is thus childless.

Note that, since  $\text{Ch}^G(v) = \emptyset$ , the marginalization  $G^{\setminus \{v\}}$  has no bi-directed edges and thus induces again a causal Bayesian network without latent variables with  $\#V^{\setminus \{v\}} = n - 1 < n$ .

Furthermore, we have the factorization:

$$P(X_V | \text{do}(X_J)) = P_v(X_v | X_{\text{Pa}^G(v)}) \otimes \underbrace{\bigotimes_{w \in \text{Pred}_{<}^G(v) \setminus J} P_w(X_w | X_{\text{Pa}^G(w)})}_{P(X_{\text{Pred}_{<}^G(v) \setminus J} | \text{do}(X_J))}.$$

This factorization implies that we already have the conditional independence:

$$X_v \perp_{P(X_V | \text{do}(X_J))} \parallel X_{\text{Pred}_{<}^G(v)} | X_D,$$

where we put  $D := \text{Pa}^G(v)$ .

In the following we will distinguish between 4 cases:

A.)  $v \in A \setminus C$ ,

B.)  $v \in B \setminus C$ ,

C.)  $v \in C$ ,

D.)  $v \notin A \cup J \cup B \cup C$ ,

Note that  $v \in V$ , thus  $v \notin J$ , which shows that the above cover all possible cases. Further note that:

$$A \perp_G^d B | C,$$

implies that:

$$A \cap (J \cup B) \subseteq C.$$

Otherwise a trivial walk from  $A$  to  $J \cup B$  would be  $C$ -open. This shows that  $A \setminus C$ ,  $(J \cup B) \setminus C$  and  $C$  are pairwise disjoint.

Case D.):  $v \notin A \cup J \cup B \cup C$ . Then we can marginalize out  $v$  and use the equivalence:

$$A \perp_G^d B | C \iff A \perp_{G \setminus v}^d B | C.$$

With  $\#V \setminus \{v\} < n$  and induction (IND) we then get:

$$X_A \perp_{P(X_V | \text{do}(X_J))} \perp\!\!\!\perp X_B | X_C.$$

This shows the claim in case D.

Case A.):  $v \in A \setminus C$ . Then we can write:

$$\begin{aligned} A &= A' \dot{\cup} (A \cap C) \dot{\cup} \{v\}, \\ B &= B' \dot{\cup} (B \cap C), \end{aligned}$$

with some disjoint  $A' \subseteq A \setminus C$  and  $B' \subseteq B \setminus C$ . We then have the implications:

$$\begin{array}{ccc} A \perp_G^d B | C & \xrightarrow{\text{Right Decomposition}} & A \perp_G^d B' | C \\ & \xrightarrow{\text{Left Decomposition}} & A' \perp_G^d B' | C \\ & \xrightarrow{\text{marginalization, } v \notin A' \cup J \cup B' \cup C} & A' \perp_{G \setminus \{v\}}^d B' | C \\ & \xrightarrow{\text{induction (IND)}} & X_{A'} \perp_{P(X_V | \text{do}(X_J))} \perp\!\!\!\perp X_{B'} | X_C. \quad (\#1) \end{array}$$

On the other hand we have with  $D = \text{Pa}^G(v)$ :

$$\begin{array}{ccc}
A \perp_G^d B \mid C & \xrightarrow{\text{Right Decomposition, } B' \subseteq B} & A \perp_G^d B' \mid C \\
& \xrightarrow{\text{Left Weak Union, } A = A' \dot{\cup} (A \cap C) \dot{\cup} \{v\}} & \{v\} \perp_G^d B' \mid A' \dot{\cup} C \\
& \xrightarrow{(*) \text{, see below}} & D \perp_G^d B' \mid A' \dot{\cup} C \\
& \xrightarrow{\text{marginalization, } v \notin D \cup J \cup B' \cup A' \cup C} & D \perp_{G \setminus \{v\}}^d B' \mid A' \dot{\cup} C \\
& \xrightarrow{\text{induction (IND)}} & X_D \perp_{P(X_V \mid \text{do}(X_J))} X_{B'} \mid X_{A' \dot{\cup} C} \\
& \xrightarrow{A' \dot{\cup} C} & X_D \perp_{P(X_V \mid \text{do}(X_J))} X_{B'} \mid X_{A'}, X_C. \quad (\#2)
\end{array}$$

(\*) holds since every  $(A' \dot{\cup} C)$ -open walk  $w \rightsquigarrow \dots$  from a  $w \in D = \text{Pa}^G(v)$  to  $J \cup B'$  extends to an  $(A' \dot{\cup} C)$ -open walk from  $v$  to  $J \cup B'$  via  $v \leftarrow w \rightsquigarrow \dots$ , as  $w$  stays a non-collider in the extended walk (not in  $A' \dot{\cup} C$ ) and  $v \notin A' \dot{\cup} C$ .

As discussed above we also already have the conditional independence:

$$X_v \perp_{P(X_V \mid \text{do}(X_J))} X_{\text{Pred}_{<}^G(v)} \mid X_D.$$

With this and  $A' \dot{\cup} B' \dot{\cup} C \subseteq \text{Pred}_{<}^G(v)$  we get the implications:

$$\begin{array}{ccc}
& & X_v \perp_{P(X_V \mid \text{do}(X_J))} X_{\text{Pred}_{<}^G(v)} \mid X_D \\
\text{Right Decomposition} & \xrightarrow{\quad} & X_v \perp_{P(X_V \mid \text{do}(X_J))} X_{A'}, X_{B'}, X_C \mid X_D \\
\text{Right Weak Union} & \xrightarrow{\quad} & X_v \perp_{P(X_V \mid \text{do}(X_J))} X_{B'} \mid X_{A'}, X_C, X_D \\
\text{Left Contraction, (\#2)} & \xrightarrow{\quad} & X_v, X_D \perp_{P(X_V \mid \text{do}(X_J))} X_{B'} \mid X_{A'}, X_C \\
\text{Left Decomposition} & \xrightarrow{\quad} & X_v \perp_{P(X_V \mid \text{do}(X_J))} X_{B'} \mid X_{A'}, X_C \\
\text{Left Contraction, (\#1)} & \xrightarrow{\quad} & X_{A'}, X_v \perp_{P(X_V \mid \text{do}(X_J))} X_{B'} \mid X_C \\
X_J\text{-Inverted Right Decomposition} & \xrightarrow{\quad} & X_{A'}, X_v \perp_{P(X_V \mid \text{do}(X_J))} X_J, X_{B'}, X_C \mid X_C \\
\text{Right Decomposition, } B \subseteq B' \dot{\cup} C & \xrightarrow{\quad} & X_{A'}, X_v \perp_{P(X_V \mid \text{do}(X_J))} X_B \mid X_C. \quad (\#3)
\end{array}$$

By (Extended) Left Redundancy we have:

$$X_{A'}, X_v, X_C \perp_{P(X_V \mid \text{do}(X_J))} X_B \mid X_{A'}, X_v, X_C.$$

With this we get the implications:

$$\begin{array}{ccc}
& & X_{A'}, X_v, X_C \perp\!\!\!\perp_{P(X_V|\text{do}(X_J))} X_B | X_{A'}, X_v, X_C \\
\text{Left Contraction, (\#3)} \longrightarrow & & \\
\text{Left Decomposition, } A \subseteq A' \dot{\cup} \{v\} \dot{\cup} C \longrightarrow & & X_{A'}, X_v, X_{A'}, X_v, X_C \perp\!\!\!\perp_{P(X_V|\text{do}(X_J))} X_B | X_C \\
& & X_A \perp\!\!\!\perp_{P(X_V|\text{do}(X_J))} X_B | X_C.
\end{array}$$

This shows the claim in case A.

Case B.):  $v \in B \setminus C$ . Then we can write:

$$\begin{aligned}
A &= A' \dot{\cup} (A \cap C), \\
B &= B' \dot{\cup} (B \cap C) \dot{\cup} \{v\},
\end{aligned}$$

with some disjoint  $A' \subseteq A \setminus C$  and  $B' \subseteq B \setminus C$ .

We then have the implications:

$$\begin{array}{ccc}
A \perp\!\!\!\perp_G^d B | C \xrightarrow{\text{Left Decomposition}} & & A' \perp\!\!\!\perp_G^d B | C \\
& \xrightarrow{\text{Right Decomposition}} & A' \perp\!\!\!\perp_G^d B' | C \\
& \xrightarrow{\text{marginalization, } v \notin A' \cup J \cup B' \cup C} & A' \perp\!\!\!\perp_{G \setminus \{v\}}^d B' | C \\
& \xrightarrow{\text{induction (IND)}} & X_{A'} \perp\!\!\!\perp_{P(X_V|\text{do}(X_J))} X_{B'} | X_C. \quad (\#1')
\end{array}$$

Again with  $D = \text{Pa}^G(v)$  we get:

$$\begin{array}{ccc}
A \perp\!\!\!\perp_G^d B | C \xrightarrow{\text{Left Decomposition}} & & A' \perp\!\!\!\perp_G^d B | C \\
& \xrightarrow{\text{Right Decomposition}} & A' \perp\!\!\!\perp_G^d B' \cup \{v\} | C \\
& \xrightarrow{\text{Right Weak Union}} & A' \perp\!\!\!\perp_G^d \{v\} | B' \dot{\cup} C \\
& \xrightarrow{(\bullet), \text{ see below}} & A' \perp\!\!\!\perp_G^d D | B' \dot{\cup} C \\
& \xrightarrow{\text{marginalization, } v \notin A' \cup J \cup D \cup B' \cup C} & A' \perp\!\!\!\perp_{G \setminus \{v\}}^d D | B' \dot{\cup} C \\
& \xrightarrow{\text{induction (IND)}} & X_{A'} \perp\!\!\!\perp_{P(X_V|\text{do}(X_J))} X_D | X_{B' \dot{\cup} C} \\
& \xrightarrow{B' \dot{\cup} C} & X_{A'} \perp\!\!\!\perp_{P(X_V|\text{do}(X_J))} X_D | X_{B'}, X_C. \quad (\#2')
\end{array}$$

(•) holds since every  $(B' \dot{\cup} C)$ -open walk  $\cdots \ast \ast w$  from  $A'$  to a  $w \in J \cup D$  extends to a  $(B' \dot{\cup} C)$ -open walk from  $A'$  to  $J \cup \{v\}$ , either because  $w \in J$  or via  $\cdots \ast \ast w \rightarrow v$  in  $w \in D = \text{Pa}^G(v)$ . Note again that  $w$  stays a non-collider in the extended walk (outside of  $B' \dot{\cup} C$ ) and  $v \notin B' \dot{\cup} C$ .

As before we will use the following conditional independence:

$$X_v \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_{\text{Pred}_{\leq}^G(v)} \mid X_D.$$

With this and  $A' \cup J \cup B' \cup C \subseteq \text{Pred}_{\leq}^G(v)$  we get the implications:

$$\begin{array}{l}
\begin{array}{l}
\text{Right Decomposition} \\
\hline\hline
\end{array} \\
\begin{array}{l}
\text{Right Weak Union} \\
\hline\hline
\end{array} \\
\begin{array}{l}
\text{Flipped Left Cross Contraction, (\#2')} \\
\hline\hline
\end{array} \\
\begin{array}{l}
\text{Right Decomposition} \\
\hline\hline
\end{array} \\
\begin{array}{l}
\text{Right Contraction, (\#1')} \\
\hline\hline
\end{array} \\
\begin{array}{l}
\text{X}_J\text{-Inverted Right Decomposition} \\
\hline\hline
\end{array} \\
\begin{array}{l}
\text{Right Decomposition, } B \subseteq B' \dot{\cup} \{v\} \dot{\cup} C \\
\hline\hline
\end{array}
\end{array}
\begin{array}{l}
X_v \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_{\text{Pred}_{\leq}^G(v)} \mid X_D \\
X_v \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_{A'}, X_{B'}, X_C \mid X_D \\
X_v \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_{A'} \mid X_{B'}, X_C, X_D \\
X_{A'} \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_D, X_v \mid X_{B'}, X_C \\
X_{A'} \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_v \mid X_{B'}, X_C \\
X_{A'} \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_{B'}, X_v \mid X_C \\
X_{A'} \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_J, X_{B'}, X_v, X_C \mid X_C \\
X_{A'} \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_B \mid X_C. \quad (\#3')
\end{array}$$

By Redundancy we have:

$$X_{A'}, X_C \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_B \mid X_{A'}, X_C.$$

With this we get the implications:

$$\begin{array}{l}
\begin{array}{l}
\text{Left Contraction, (\#3')} \\
\hline\hline
\end{array} \\
\begin{array}{l}
\text{Left Decomposition, } A \subseteq A' \dot{\cup} C \\
\hline\hline
\end{array}
\end{array}
\begin{array}{l}
X_{A'}, X_C \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_B \mid X_{A'}, X_C \\
X_{A'}, X_{A'}, X_C \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_B \mid X_C. \\
X_A \perp\!\!\!\perp_{P(X_V | \text{do}(X_J))} X_B \mid X_C.
\end{array}$$

This shows the claim in case B.

Case C.):  $v \in C$ . Then we can write:

$$\begin{aligned}
A &= A' \dot{\cup} (A \cap C), \\
B &= B' \dot{\cup} (B \cap C), \\
C &= C' \dot{\cup} \{v\},
\end{aligned}$$

with some pairwise disjoint  $A' \subseteq A \setminus C$ ,  $B' \subseteq B \setminus C$  and  $C' \subseteq C$ .

We then get the implications.

$$\begin{array}{ccc}
A \underset{G}{\perp}^d B | C & \xrightarrow{\text{Left Decomposition}} & A' \underset{G}{\perp}^d B | C \\
& \xrightarrow{\text{Right Decomposition}} & A' \underset{G}{\perp}^d B' | C \\
& \xrightarrow{C=C' \dot{\cup} \{v\}} & A' \underset{G}{\perp}^d B' | C' \dot{\cup} \{v\}
\end{array}$$

We now claim that:

$$A' \underset{G}{\perp}^d B' | C' \dot{\cup} \{v\}$$

implies that one of the following statements holds:

$$A' \dot{\cup} \{v\} \underset{G}{\perp}^d B' | C' \quad \vee \quad A' \underset{G}{\perp}^d B' \dot{\cup} \{v\} | C'.$$

Assume the contrary:

$$A' \dot{\cup} \{v\} \not\underset{G}{\perp}^d B' | C' \quad \wedge \quad A' \not\underset{G}{\perp}^d B' \dot{\cup} \{v\} | C'.$$

So there exist shortest  $C'$ -open walks  $\pi_1$  and  $\pi_2$  in  $G$  such that all colliders are in  $C'$ :

$$\pi_1 : \quad A' \cup \{v\} \ni u_0 \rightsquigarrow \dots \rightsquigarrow u_k \in J \cup B',$$

and:

$$\pi_2 : \quad A' \ni w_0 \rightsquigarrow \dots \rightsquigarrow w_m \in J \cup (B' \dot{\cup} \{v\}).$$

So all non-colliders of  $\pi_1$  and  $\pi_2$  are outside of  $C'$ . Since we consider shortest walks and  $v \notin C'$  at most an end node of  $\pi_1$  and  $\pi_2$  could be equal to  $v$ . Otherwise one could shorten the walk.

Then note that  $v \notin A'$  and  $v \notin J \cup B'$ , thus:  $u_k \neq v$  and  $w_0 \neq v$ .

If now  $\pi_i$  does not contain  $v$  as an (end) node, then  $\pi_i$  would be  $(C' \dot{\cup} \{v\})$ -open, which is a contradiction to the assumption:

$$A' \underset{G}{\perp}^d B' | C' \dot{\cup} \{v\}.$$

So we can assume that the other end nodes equal  $v$ , i.e.:  $u_0 = v$  and  $w_m = v$ .

Furthermore, both  $\pi_1$  and  $\pi_2$  are non-trivial walks, since  $u_0 \neq u_k$  and  $w_0 \neq w_m$ . Since  $v$  is childless and  $k, m \geq 1$  we have that the  $\pi_i$  are of the forms:

$$\pi_1 : \quad v \longleftarrow u_1 \rightsquigarrow \dots \rightsquigarrow u_k,$$

and:

$$\pi_2 : \quad w_0 \rightsquigarrow \dots \rightsquigarrow w_{m-1} \longrightarrow v,$$

with  $u_1, w_{m-1} \in D = \text{Pa}^G(v)$ . Then the following walk:

$$A' \ni w_0 \ast\ast \dots \ast\ast w_{m-1} \rightarrow v \leftarrow u_1 \ast\ast \dots \ast\ast u_k \in J \cup B',$$

is a  $(C' \dot{\cup} \{v})$ -open walk from  $A'$  to  $J \cup B'$ , in contradiction to:

$$A' \perp_G^d B' \mid C' \dot{\cup} \{v\}.$$

So the claim:

$$A' \dot{\cup} \{v\} \perp_G^d B' \mid C' \quad \vee \quad A' \perp_G^d B' \dot{\cup} \{v\} \mid C',$$

must be true. So we reduced case C to case A or case B, which then imply:

$$X_A, X_v \perp_{P(X_V \mid \text{do}(X_J))} X_B \mid X_{C'} \quad \vee \quad X_A \perp_{P(X_V \mid \text{do}(X_J))} X_B, X_v \mid X_{C'}.$$

If we apply Left Weak Union to the left and Right Weak Union to the right we get:

$$X_A \perp_{P(X_V \mid \text{do}(X_J))} X_B \mid X_{C'}, X_v,$$

which implies:

$$X_A \perp_{P(X_V \mid \text{do}(X_J))} X_B \mid X_{C'}.$$

This shows the claim in case C. □

### 4.3. Operations on Causal Bayesian Networks

#### 4.3.1. Hard Interventions on Causal Bayesian Networks

**Definition 4.3.1** (Hard intervention on causal Bayesian network). *Consider a causal Bayesian network (CBN) given by:*

$$M = \left( G = (J, V, E), (P_v(X_v \mid X_{\text{Pa}^G(v)}))_{v \in V} \right).$$

Now let  $W \subseteq J \cup V$  be any subset. Then we define the intervened causal Bayesian network w.r.t.  $W$  via:

1. CDAG:  $G_{\text{do}(W)} = (J \cup W, V \setminus W, E_{\text{do}(W)})$ , and:
2. Markov kernels:  $P_v(X_v \mid X_{\text{Pa}^G(v)})$  for  $v \in V \setminus W$ .

Its observable Markov kernel is then:

$$P(X_{V \setminus W} \mid \text{do}(X_{J \cup W})) = \bigotimes_{v \in V \setminus W} P_v(X_v \mid X_{\text{Pa}^G(v)}).$$

Note that if  $v \in V \setminus W$  then  $\text{Pa}^G(v) = \text{Pa}^{G_{\text{do}(W)}}(v)$ .

**Remark 4.3.2.** Note that the above notations imply for every  $v \in V$  and  $W \subseteq V \setminus \{v\}$  the identifications:

$$P_v(X_v|X_{\text{Pa}^G(v)}) = P(X_v|\text{do}(X_{J \cup V \setminus W})) = P(X_v|\text{do}(X_{\text{Pa}^G(v)})),$$

which we will use interchangeably in the following.

**Remark 4.3.3** (Hard intervention on causal Bayesian network with latent variables). We define hard interventions on an L-CBN the same way as on a CBN, but we usually only allow for interventions on sets  $W \subseteq J \cup V$ , i.e. with  $W \cap U = \emptyset$ , where  $U$  is the set of latent variables.

### 4.3.2. Node-Splitting Hard Interventions on Causal Bayesian Networks

**Definition 4.3.4** (Node-splitting hard intervention on causal Bayesian network). Consider a causal Bayesian network (CBN) given by  $(G, P(X_V|\text{do}(X_J)))$  with CDAG:  $G = (J, V, E)$  and Markov kernels:  $P_v(X_v|X_{\text{Pa}^G(v)})$  for  $v \in V$ . Now let  $W \subseteq V$  be any subset. Then we define the node-splitting hard intervention w.r.t.  $W$  as the causal Bayesian network given by:

1. CDAG:  $G' := G_{\text{swig}(W)} = (J \dot{\cup} W^i, W^o \dot{\cup} V \setminus W, E_{\text{swig}(W)})$ , and:
2. Markov kernels for  $v \in V$ :

$$P_{v^o}(X_{v^o} \in A | X_{\text{Pa}^{G'}(v^o)} = \tilde{x}) := P_v(X_v \in A | X_{\text{Pa}^G(v)} = \tilde{x}),$$

where for brevity we put  $v^o := v$  for  $v \in V \setminus W$ .

**Remark 4.3.5.** Similarly, we can define node-splitting hard interventions on causal Bayesian network with latent variables, but allow only  $W$  with  $W \cap U = \emptyset$ .

### 4.3.3. Soft Interventions on Causal Bayesian Networks

**Remark 4.3.6** (Modelling soft interventions on causal Bayesian networks). Consider a causal Bayesian network given by  $(G, P(X_V|\text{do}(X_J)))$  with CDAG:  $G = (J, V, E)$  and Markov kernels:  $P_v(X_v|X_{\text{Pa}^G(v)})$  for  $v \in V$ .

Let  $W \subseteq J \cup V$ . In order to model a soft intervention on variables  $X_w$  for  $w \in W \setminus J$ , we introduce intervention nodes  $I_w \rightarrow w$  for  $w \in W \setminus J$ , which come with new input variables  $X_{I_w}$ , and replace the Markov kernel:

$$P_w(X_w|X_{\text{Pa}^G(w)})$$

for  $w \in W \setminus J$  by one that models the dependence on the soft intervention variables properly:

$$P_w(X_w|X_{\text{Pa}^G(w)}, X_{I_w}).$$

So the softly intervened causal Bayesian network w.r.t.  $W$  then has:



1. CDAG:  $G_{\text{do}(I_W)} = (J \dot{\cup} \{I_w \mid w \in W\}, V, E \dot{\cup} \{I_w \rightarrow w \mid w \in W \setminus J\})$ , and:
2. Markov kernels:

$$P_v(X_v \mid X_{\text{Pa}^G(v)}) \text{ for } v \in V \setminus W, \text{ and:}$$

$$P_w(X_w \mid X_{\text{Pa}^G(w)}, X_{I_w}) \text{ for } w \in W \setminus J.$$

Note that  $\text{Pa}^{G_{\text{do}(I_W)}}(w) = \text{Pa}^G(w) \dot{\cup} \{I_w\}$  for  $w \in W \setminus J$  and  $\text{Pa}^{G_{\text{do}(I_W)}}(v) = \text{Pa}^G(v)$  for  $v \in V \setminus W$ .

**Remark 4.3.7** (Modelling hard interventions with intervention nodes). *It is sometimes beneficial to model hard interventions with intervention nodes. Let the setting be like in Remark 4.3.6. When we model hard interventions with intervention nodes we make the further more specific choices for  $w \in W \setminus J$ :*

1.  $\mathcal{X}_{I_w} := \mathcal{X}_w \dot{\cup} \{\star\}$ ,
2.  $P_w(X_w \in A \mid X_{\text{Pa}^G(w)} = x_{\text{Pa}^G(w)}, X_{I_w} = x_{I_w}) :=$ 

$$\begin{cases} P_w(X_w \in A \mid X_{\text{Pa}^G(w)} = x_{\text{Pa}^G(w)}), & \text{if } x_{I_w} = \star, \\ \delta(X_w \in A \mid X_w = x_{I_w}) = \mathbb{1}_A(x_{I_w}), & \text{if } x_{I_w} \neq \star. \end{cases}$$

Note that the CDAG will then rather be:  $G_{\text{do}(I_W)}$  in contrast to:  $G_{\text{do}(W)}$ .

The above choices reflect that if we put  $X_{I_w} = \star$  then no intervention occurs and the value of  $X_w$  is (probabilistically) determined using the usual Markov kernel. But if we put  $X_{I_w} = x_{I_w} \neq \star$  then we change the value of  $X_w$  to  $x_{I_w}$  (with 100% probability) independent of the values of its parents. This is then similar to the hard intervention:  $\text{do}(X_w = x_{I_w})$ . This allows us to model simultaneously the unintervened and an intervened version of the CBN with a single CBN.

**Remark 4.3.8.** *Again, we can do all the above also with causal Bayesian network with latent variables, but allow only  $W$  with  $W \cap U = \emptyset$ .*

#### 4.3.4. Marginalization of Causal Bayesian Networks

**Definition 4.3.9** (Marginalization of causal Bayesian network with latent variables). *Consider a causal Bayesian network with latent variables (L-CBN):*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v \mid X_{\text{Pa}^{G^+}(v)}) \right)_{v \in V \cup U} \right),$$

Let  $W \subseteq V$  be a subset. We then define the marginalized L-CBN by just replacing  $V$  with  $V \setminus W$  and  $U$  with  $U \dot{\cup} W$ . The Markov kernels  $P_v$  for  $v \in V \dot{\cup} U = (V \setminus W) \dot{\cup} (U \dot{\cup} W)$  stay the same.

With this definition the observable Markov kernel marginalizes to:

$$P(X_{V \setminus W} \mid \text{do}(X_J)),$$

and the observable CADMG becomes:

$$(G^+) \setminus (U \dot{\cup} W) = G \setminus W,$$

i.e. the marginalized  $G$  w.r.t.  $W$ .

## 4.4. Standard Forms of Causal Bayesian Networks

**Definition 4.4.1.** Consider two causal Bayesian network with latent variables (L-CBNs):

$$M_1 = \left( G_1^+ = (J_1, (V_1, U_1), E_1^+), \left( P_{1,v}(X_v | X_{\text{Pa}G_1^+(v)}) \right)_{v \in U_1 \cup V_1} \right),$$

$$M_2 = \left( G_2^+ = (J_2, (V_2, U_2), E_2^+), \left( P_{2,v}(X_v | X_{\text{Pa}G_2^+(v)}) \right)_{v \in U_2 \cup V_2} \right).$$

We call them interventionally equivalent if all of the following conditions hold:

1.  $J_1 = J_2 =: J$ ,
2.  $V_1 = V_2 =: V$ ,
3.  $\mathcal{X}_{1,v} = \mathcal{X}_{2,v} =: \mathcal{X}_v$  for all  $v \in J \cup V$ ,
4. for all subsets  $W \subseteq V$  we have the equality of the intervened Markov kernels:

$$P_1(X_{V \setminus W} | \text{do}(X_{J \cup W})) = P_2(X_{V \setminus W} | \text{do}(X_{J \cup W})).$$

**Definition 4.4.2** (Cliques and maximal cliques of undirected graphs). Let  $G = (V, L)$  be an undirected graph. A set of nodes  $W \subseteq V$  is called a clique<sup>16</sup> of  $G$  if for all  $w_1, w_2 \in W$  with  $w_1 \neq w_2$  we have that the edge  $w_1 - w_2 \in L$ . A clique  $W$  is called a maximal clique of  $G$  if for every clique  $\tilde{W}$  of  $G$  with  $W \subseteq \tilde{W}$  we have that  $W = \tilde{W}$ .

**Definition/Theorem 4.4.3** (Standard forms of L-CBNs). Consider a causal Bayesian network  $M$  with latent variables (L-CBN) with observable CADMG  $G = (J, V, E, L)$ . Let

$$\mathcal{C} := \{W \subseteq V \mid W \text{ maximal clique of } (V, L)\}.$$

the sets of all maximal cliques of the (undirected) graph consisting only of the nodes from  $V$  and the bi-directed edges from  $G$ . Define the set of (latent) nodes:

$$\tilde{U} := \{\tilde{u}_W \mid W \in \mathcal{C}\},$$

and directed edges:

$$\tilde{E}^+ := E \cup \{\tilde{u}_W \rightarrow w \mid W \in \mathcal{C}, w \in W\}.$$

Then there exists an L-CBN of the form:

$$\tilde{M} = \left( \tilde{G}^+ = (J, (V, \tilde{U}), \tilde{E}^+), \left( \tilde{P}_v(X_v | X_{\text{Pa}\tilde{G}^+(v)}) \right)_{v \in V \cup \tilde{U}} \right)$$

that is interventionally equivalent to  $M$ . Furthermore, we can choose to arrange them in one of the following ways:

---

<sup>16</sup>A clique is also called *complete subgraph* in the literature.

1. Structural causal model form: All Markov kernels for  $v \in V$  are deterministic:

$$\tilde{P}_v(X_v \in A | X_{\text{Pa}^{G^+}(v)} = \tilde{x}) = \delta(R_v \in A | X_{\text{Pa}^{G^+}(v)} = \tilde{x}),$$

for some measurable maps  $R_v$ ,  $v \in V$ . OR:

2. Canonical form: All latent variables  $\tilde{u}_W$  with  $\#W = 1$  and the corresponding variables, edges and Markov kernels can be removed from  $\tilde{M}$  as well, leaving us only with the latent variables  $\tilde{u}_W$  with  $W \in \mathcal{C}$  and  $\#W \geq 2$ .

**Remark 4.4.4.** Consider the standard forms  $\tilde{M}$  of  $M$  from Definition/Theorem 4.4.3.

1. In particular, we have:

- a)  $(\tilde{G}^+) \setminus \tilde{U} = G$ ,

- b)  $\text{Pa}^{\tilde{G}^+}(\tilde{U}) = \emptyset$ ,

- c)  $\text{Ch}^{\tilde{G}^+}(u) \in \mathcal{C}$  for every  $u \in \tilde{U}$ .

2. We can use measurable embeddings/isomorphisms:  $\mathcal{X}_u \hookrightarrow [0, 1]$  for  $u \in \tilde{U}$  to further restrict to the case:

- a)  $\mathcal{X}_u \cong [0, 1]$ ,

- b)  $\tilde{P}_u(X_u)$  is the uniform distribution on  $[0, 1]$ .

3. Note that the Markov kernels dependent on  $X_{\tilde{U}}$  might not be unique as we can always transform  $[0, 1]$  to  $[0, 1]$  in strange ways.
4. The construction of the canonical form generally<sup>17</sup> leads to an interventionally equivalent L-CBN with the smallest number of latent variables such that its observable CADMG stays unchanged.
5. The construction of the structural causal model form generally<sup>17</sup> leads to an interventionally equivalent L-CBN with the smallest number of latent variables such that its observable CADMG stays unchanged and such that every Markov kernel with non-trivial input is deterministic.

**Remark 4.4.5** (Marginalizations and hard interventions on standard forms). Let the following L-CBN be in one of the standard forms:

$$(G^+ = (J, V, U, E^+), P(X_{V \cup U} | \text{do}(X_J))).$$

Now let  $W \subseteq V$  then we defined the marginalization w.r.t.  $W$  by replacing  $V$  with  $V \setminus W$  and  $U$  with  $U \cup W$ . We could re-define the marginalization as a corresponding standard form of that procedure.

Similarly we could post-process hard interventions with standardization steps.

---

<sup>17</sup>Excluding degenerate L-CBNs. In those cases one could possibly remove even more latent variables.

## Proofs - Standard Forms of Causal Bayesian Networks

*Proof.* Step 1. For every  $v \in V \cup U$  we can write the Markov kernel  $P_v$  as the composition of a deterministic one and a uniform distribution  $\bar{P}_{\bar{v}}(X_{\bar{v}})$  on  $\mathcal{X}_{\bar{v}} := [0, 1]$  by Remark 2.7.4:

$$P_v(X_v | X_{\text{Pa}^{G^+}(v)}) = \delta(R_v | X_{\bar{v}}, X_{\text{Pa}^{G^+}(v)}) \circ \bar{P}_{\bar{v}}(X_{\bar{v}}).$$

We now put:

$$\bar{U} := U \dot{\cup} \{\bar{v} \mid v \in V \cup U\}, \quad \bar{E}^+ := E^+ \dot{\cup} \{\bar{v} \rightarrow v \mid v \in V \cup U\},$$

and to get  $\bar{M}$  we add the  $\bar{P}_{\bar{v}}$  to  $M$  and replace  $P_v$  for  $v \in V \cup U$  by the deterministic one given by:

$$\bar{P}_v(X_v \in A | X_{\bar{v}}, X_{\text{Pa}^{G^+}(v)}) := \delta(R_v \in A | X_{\bar{v}}, X_{\text{Pa}^{G^+}(v)}).$$

Then  $\bar{G}^+$  clearly marginalizes to  $G^+$  (when we marginalize out all the  $\bar{v}$  again) and the marginal of:

$$\bar{P}_v(X_v \in A | X_{\bar{v}}, X_{\text{Pa}^{G^+}(v)}) \otimes \bar{P}_{\bar{v}}(X_{\bar{v}}),$$

in the defining product of the joint Markov kernel is  $P_v(X_v | X_{\text{Pa}^{G^+}(v)})$  for all  $v \in V \cup U$  by construction again.

Step 2. Marginalize out all  $u \in U$ . Let us first look at the Markov kernel side if we marginalize out  $X_u$  in the defining product of the joint Markov kernel for  $u \in U$ :

$$\begin{aligned} & \int_{\mathcal{X}_u} \bigotimes_{v \in \text{Ch}^{\bar{G}^+}(u)} \bar{P}_v(X_v | X_{\text{Pa}^{\bar{G}^+}(v) \setminus \{u\}}, X_u = x_u) \delta(R_u \in dx_u | X_{\text{Pa}^{\bar{G}^+}(u)}) \\ &= \bigotimes_{v \in \text{Ch}^{\bar{G}^+}(u)} \bar{P}_v(X_v | X_{\text{Pa}^{\bar{G}^+}(v) \setminus \{u\}}, X_u = R_u(X_{\text{Pa}^{\bar{G}^+}(u)})), \end{aligned}$$

which is again a product (only) because we marginalized a deterministic Markov kernel out. So we define:

$$\hat{P}_v(X_v | X_{\text{Pa}^{\hat{G}^+}(v)}) := \bar{P}_v(X_v | X_{\text{Pa}^{\bar{G}^+}(v) \setminus \{u\}}, X_u = R_u(X_{\text{Pa}^{\bar{G}^+}(u)})),$$

which is as the composition of deterministic Markov kernels again a deterministic Markov kernel. From this we also read off that we need to consider the graph  $\hat{G}^+$  with:

$$\text{Pa}^{\hat{G}^+}(v) := \text{Pa}^{\bar{G}^+}(v) \setminus \{u\} \cup \text{Pa}^{\bar{G}^+}(u),$$

i.e. the CDAG from  $(\bar{G}^+) \setminus U$  where we removed all bi-directed edges, and with latent nodes  $\hat{U} = \bar{U} \setminus U$ . Then note that for  $u \in U$  we have:

$$\text{Ch}^{\hat{G}^+}(u) = \text{Ch}^{(\bar{G}^+) \setminus U}(\bar{u}).$$

This implies that we recover the removed bi-directed edges if we further marginalize out all the  $\bar{u}$ , i.e.:

$$(\hat{G}^+) \setminus \hat{U} = (\bar{G}^+) \setminus \bar{U} = (G^+) \setminus U = G.$$

Step 3. We marginalize out all nodes  $u \in \hat{U}$  with  $\text{Ch}^{\hat{G}^+}(u) = \emptyset$ . For those  $u$  we have:

$$\hat{P}(X_V, X_{\hat{U}} | \text{do}(X_J)) = \hat{P}_u(X_u | X_{\text{Pa}^{\hat{G}^+}(u)}) \otimes \hat{P}(X_V, X_{\hat{U} \setminus \{u\}} | \text{do}(X_J)).$$

So marginalizing out  $X_u$  does not interfere with the rest of the Markov kernels. So from now on we can w.l.o.g. assume that  $\#\text{Ch}^{\hat{G}^+}(u) \geq 1$  for all  $u \in \hat{U}$ .

Step 4. Since  $(\hat{G}^+) \setminus \hat{U} = G$  we have that for each  $u \in \hat{U}$  the set  $\text{Ch}^{\hat{G}^+}(u)$  is a clique of  $(V, L)$ . So we can (arbitrarily) assign  $u$  to any maximal clique  $W$  of  $(V, L)$  with  $\text{Ch}^{\hat{G}^+}(u) \subseteq W$ . So let  $W$  be a fixed maximal clique of  $(V, L)$  and  $u_1, \dots, u_k \in \hat{U}$  be all  $u \in \hat{U}$  that we assigned to  $W$ . Then we consider the space:

$$\mathcal{X}_{\tilde{u}_W} := \prod_{\ell=1}^k \mathcal{X}_{u_\ell},$$

and the variables:

$$X_{\tilde{u}_W} := (X_{u_\ell})_{\ell=1, \dots, k}.$$

Then every Markov kernel dependent on such an  $X_{u_\ell}$  can be written as a Markov kernel dependent on  $X_{\tilde{u}_W}$ , by only using the  $u_\ell$  component. We will then replace  $u_1, \dots, u_k$  by the single node  $\tilde{u}_W$  and every edge of form  $u_\ell \rightarrow v$  by  $\tilde{u}_W \rightarrow v$ . If we do this for all  $u \in \hat{U}$  and maximal cliques  $W$  of  $(V, L)$  we arrive at the CADMG  $\tilde{G}^+ = (J, V, \tilde{U}, \tilde{E}^+)$ , with:

$$\tilde{E}^+ := E \cup \{\tilde{u}_W \rightarrow w \mid W \in \mathcal{C}, w \in W\}.$$

So we arrived at the desired structural causal model form and one can convince oneself that at each step we get an interventionally equivalent L-CBN to the step before.

The canonical form follows from the structural causal model form by marginalizing out all  $X_u$  with  $\#\text{Ch}^{\tilde{G}^+}(u) \leq 1$ , i.e. by replacing the left (deterministic) Markov kernel dependent on  $X_u$  in the product:

$$P_v(X_v | X_{\text{Pa}^{\tilde{G}^+}(v)} \setminus \{u\}, X_u) \otimes P_u(X_u),$$

by the composition:

$$P_v(X_v | X_{\text{Pa}^{\tilde{G}^+}(v)} \setminus \{u\}, X_u) \circ P_u(X_u),$$

which then might not be deterministic anymore. □

## 5. Identification of Causal Effects in CBNs

This section investigates under which circumstances one can *identify causal effects* and estimate them just from observational data alone under the (strong) assumption that the underlying causal graph is known. More generally, we ask the question when an interventional Markov kernel of a causal Bayesian network can be identified from the causal graph  $G$  and the observational Markov kernel alone.

We will see that the main tool to allow for such statements is the global Markov property, see Theorem 4.2.1, applied to the causal Bayesian network that is augmented with further intervention variables.

We first study under which graphical conditions interventions don't have an effect or when one essentially can replace interventions with conditioning operations. These rules will be summarized as the *three rules of do-calculus*. The main references are [Pea93a, Pea93b, Pea09], also see [Pea95, FM20, For21].

We then study under which graphical criteria one gets explicit *adjustment formulas* to estimate interventional Markov kernels from observational ones. The literature mentions the *backdoor criterion*, see [Pea93a, Pea93b, Pea09], the *extended backdoor criterion*, see [PP10, SdWR10], the *selection backdoor criterion*, see [BTP14], criteria for *selection without/partial external data*, see [CB17, CTB18], and all their generalizations to the cyclic case, see [FM20], also see [SP06a, PTKM15, For21].

Finally, we present the *ID-algorithm*, which can decide just by processing the causal graph  $G$  if an interventional Markov kernel can be identified by the observational one (under further assumptions, like strict positivity, etc.). If the algorithm does not output FAIL then it also presents a formula to estimate the queried interventional Markov kernel. The main references for the ID-algorithm are [Pea09, GP95, Tia02, TP02, Tia04, SP06b, HV06, HV08, RERS23, FM20].

### 5.1. Do-Calculus

**Remark 5.1.1** (Recap). *Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}^{G^+}(v)})) \right)_{v \in V \cup U} \right).$$

*Then we get the joint Markov kernel over all input, observed and unobserved output variables as follows:*

$$P(X_V, X_U, X_J | \text{do}(X_J)) := \bigotimes_{v \in U \cup V} P_v \left( X_v | \text{do}(X_{\text{Pa}^{G^+}(v)}) \right) \otimes \bigotimes_{j \in J} \delta(X_j | X_j).$$

*Further, for  $D \subseteq J \cup V$  and  $C \subseteq V \setminus D$  we get the combined hard and soft interventions:*

$$P(X_{V \setminus D}, X_U, X_{J \cup D}, X_{I_C} | \text{do}(X_{I_C}, X_{J \cup D})) := \bigotimes_{v \in V \setminus (C \cup D)} P_v \left( X_v | \text{do}(X_{\text{Pa}^{G^+}(v)}) \right) \otimes \bigotimes_{v \in C} P_v \left( X_v | \text{do}(X_{\text{Pa}^{G^+}(v)}, X_{I_v}) \right) \otimes$$

$$\bigotimes_{v \in U} P_v \left( X_v \mid \text{do}(X_{\text{Pa}^{G^+}(v)}) \right) \otimes \bigotimes_{j \in J \cup D} \delta(X_j \mid X_j) \otimes \bigotimes_{v \in C} \delta(X_{I_v} \mid X_{I_v}),$$

where we need to reorder all the factors such that the product is in reverse order of a topological order and where we use the following Markov kernels to model hard interventions as soft interventions,  $v \in C$ :

$$P_v \left( X_v \mid \text{do} \left( X_{\text{Pa}^{G^+}(v)}, X_{I_v} = x_{I_v} \right) \right) := \begin{cases} P_v \left( X_v \mid \text{do} \left( X_{\text{Pa}^{G^+}(v)} \right) \right), & \text{if } x_{I_v} = \star, \\ \delta(X_v \mid X_v = x_{I_v}), & \text{if } x_{I_v} \neq \star. \end{cases}$$

Finally we can also marginalize (i.e. integrating out) and condition to get:

$$P(X_A \mid X_B, \text{do}(X_{J \cup D}, X_{I_C})),$$

for any  $A, B, C, D \subseteq J \cup V$ .

For more suggestive formulas later on we also freely permute the order of symbols behind the conditioning line, e.g.:

$$P(X_A \mid \text{do}(X_F), X_B, \text{do}(X_D)) := P(X_A \mid X_B, \text{do}(X_D, X_F)).$$

Please note that no matter in which order we write the do-part and conditioning part behind the conditioning line  $\mid$ , we always assume that we perform the intervention (do) first and afterwards condition.

We will also make use of the following CADMG:

$$G_{\text{do}(I_C, D)} = (G_{\text{do}(I_C, D)}^+)^{\setminus U}.$$

W.l.o.g. we can assume:  $C \cap D = \emptyset$ .

**Theorem 5.1.2.** [Almost-sure do-calculus—in detail] Consider an L-CBN:

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v \mid \text{do}(X_{\text{Pa}^{G^+}(v)})) \right)_{v \in V \cup U} \right).$$

Let  $A, B, C \subseteq V$  and  $D \subseteq J \cup V$  be such that  $A, B, C, D$  are pairwise disjoint.

Further assume that we have reference measures  $\mu_v$  on  $\mathcal{X}_v$  for every  $v \in V$  that are each equivalent to a probability measure (in terms of absolute continuity).<sup>18</sup> We then put  $\mu_F := \bigotimes_{v \in F} \mu_v$  for  $F \subseteq V$ .

1. Insertion/deletion of observation: Assume:

$$A \underset{G_{\text{do}(D)}}{\perp^d} B \mid C \cup D.$$

<sup>18</sup>Recall the connection between absolute continuity and strictly positive densities in Corollary 2.3.20. All  $\sigma$ -finite measures satisfy this assumption.

For a fixed finite index set  $I$  consider subsets  $B^{(i)} \subseteq B$ , for  $i \in I$ , and pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A | X_{B^{(i)}}, X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{B \cup C \cup D \cup J} \rightarrow \mathcal{X}_{B^{(i)} \cup C \cup D \cup J} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_{B^{(i)}}, X_C | \text{do}(X_{D \cup J}))$ . Then there exists a measurable  $P(X_B, X_C | \text{do}(X_{D \cup J}))$ -null set  $N \subseteq \mathcal{X}_{B \cup C \cup D \cup J}$ , such that all those Markov kernels are equal on the complement  $N^c$ .

Note that if  $\mu_{B \cup C} \ll P(X_B, X_C | \text{do}(X_{D \cup J}))$  then  $N$  is also a  $\mu_{B \cup C}$ -null set, i.e. for every  $x_{D \cup J} \in \mathcal{X}_{D \cup J}$  we have:  $\mu_{B \cup C}(N_{x_{D \cup J}}) = 0$ .

If we also have the reverse  $P(X_B, X_C | \text{do}(X_{D \cup J})) \ll \mu_{B \cup C}$  then we can change the above conditional Markov kernels on a  $\mu_{B \cup C}$ -null set  $N$  while they remain versions of the corresponding conditional Markov kernel.<sup>19</sup>

2. Action/observation exchange: Assume:

$$A \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B | B \cup C \cup D.$$

For a fixed finite index set  $I$  consider decompositions  $B = B_1^{(i)} \dot{\cup} B_2^{(i)}$ , for  $i \in I$ , and pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A | X_{B_1^{(i)}}, \text{do}(X_{B_2^{(i)}}), X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{B \cup C \cup D \cup J} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \otimes \mu_{B_2^{(i)}}$  and assume the following absolute continuities:

$$\mu_{B \cup C} \ll P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \otimes \mu_{B_2^{(i)}}$$

for all  $i \in I$ .<sup>20</sup> Then there exists a measurable  $\mu_{B \cup C}$ -null set  $N \subseteq \mathcal{X}_{B \cup C \cup D \cup J}$ , such that all those conditional Markov kernels are equal on the complement  $N^c$ .

If we also assume the reverse absolute continuities for all  $i \in I$ :

$$P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \otimes \mu_{B_2^{(i)}} \ll \mu_{B \cup C},$$

<sup>19</sup>Note that the absolute continuities:  $\mu_{B \cup C} \ll P(X_B, X_C | \text{do}(X_{D \cup J})) \ll \mu_{B \cup C}$  hold if  $P(X_B, X_C | \text{do}(X_{D \cup J}))$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $\mu_{B \cup C}$ . Furthermore, the converse is also true for  $\sigma$ -finite reference measures  $\mu_{B \cup C}$  by Corollary 2.3.20.

<sup>20</sup>If you instead expected to pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A | X_{B_1^{(i)}}, \text{do}(X_{B_2^{(i)}}), X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{B \cup C \cup D \cup J} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}), X_{D \cup J})$  and to assume the absolute continuities

$$\mu_{B_1^{(i)} \cup C} \ll P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \ll \mu_{B_1^{(i)} \cup C}$$

for all  $i \in I$ : that would lead to a similar, but slightly weaker statement.



then all those conditional Markov kernels are versions of each other.<sup>21</sup>

3. Insertion/deletion of action: Assume:

$$A \underset{G_{\text{do}(I_B, D)}}{\overset{d}{\perp}} I_B \mid C \cup D.$$

For a fixed finite index set  $I$  consider subsets  $B^{(i)} \subseteq B$ , for  $i \in I$ , and pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A \mid \text{do}(X_{B^{(i)}}), X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{B \cup C \cup D \cup J} \rightarrow \mathcal{X}_{B^{(i)} \cup C \cup D \cup J} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_C \mid \text{do}(X_{B^{(i)}}), X_{D \cup J})$  and assume the following absolute continuities:

$$\mu_C \ll P(X_C \mid \text{do}(X_{B^{(i)}}), X_{D \cup J})$$

for all  $i \in I$ . Then there exists a measurable  $\mu_C$ -null set  $N \subseteq \mathcal{X}_{B \cup C \cup D \cup J}$ , such that all those conditional Markov kernels are equal on the complement  $N^c$ .

If we also assume the reverse absolute continuities for all  $i \in I$ :

$$P(X_C \mid \text{do}(X_{B^{(i)}}), X_{D \cup J}) \ll \mu_C,$$

then all those conditional Markov kernels are versions of each other.<sup>22</sup>

The proof can be found in at the end of this section.

We now summarize on how to apply Theorem 5.1.2 more concretely as a corollary.

**Corollary 5.1.3** (Almost-sure do-calculus—simplified). *Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v \mid \text{do}(X_{\text{Pa}G^+(v)})) \right)_{v \in V \cup U} \right).$$

Further assume that we have reference measures  $\mu_v$  on  $\mathcal{X}_v$  for every  $v \in V$ . We then put  $\mu_F := \bigotimes_{v \in F} \mu_v$  for  $F \subseteq V$ . Let  $A, B, C \subseteq V$  and  $D \subseteq J \cup V$  be such that  $A, B, C, D$  are pairwise disjoint. Then we have the following 3 rules relating marginal conditional to marginal interventional Markov kernels:

<sup>21</sup>Note that the absolute continuities:  $\mu_{B \cup C} \ll P(X_{B_1^{(i)}}, X_C \mid \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \otimes \mu_{B_2^{(i)}} \ll \mu_{B \cup C}$  hold if the absolute continuities:  $\mu_{B_1^{(i)} \cup C} \ll P(X_{B_1^{(i)}}, X_C \mid \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \ll \mu_{B_1^{(i)} \cup C}$  hold, which hold if  $P(X_{B_1^{(i)}}, X_C \mid \text{do}(X_{B_2^{(i)}}), X_{D \cup J})$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $\mu_{B_1^{(i)} \cup C}$ . Furthermore, the converse is also true for  $\sigma$ -finite reference measures  $\mu_{B_1^{(i)} \cup C}$  by Corollary 2.3.20.

<sup>22</sup>Note that absolute continuities:  $\mu_C \ll P(X_C \mid \text{do}(X_{B^{(i)}}), X_{D \cup J}) \ll \mu_C$  hold if  $P(X_C \mid \text{do}(X_{B^{(i)}}), X_{D \cup J})$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $\mu_C$ . Furthermore, the converse is also true for  $\sigma$ -finite reference measures  $\mu_C$  by Corollary 2.3.20.

1. Insertion/deletion of observation, for  $J \subseteq D$ : Assume that we want to establish the a.s.-equality:

$$P(X_A|X_B, X_C, \text{do}(X_D)) = P(X_A|X_C, \text{do}(X_D)) \quad \mu_{B \cup C}\text{-a.s.},$$

then it is sufficient to assume/check the following  $d$ -separation and absolute continuities:

$$A \underset{G_{\text{do}(D)}}{\perp^d} B | C \cup D, \quad \mu_{B \cup C} \ll P(X_B, X_C | \text{do}(X_D)) \ll \mu_{B \cup C}.$$

2. Action/observation exchange, for  $J \subseteq D$ : Assume that we want to establish the a.s.-equality:

$$P(X_A|X_B, X_C, \text{do}(X_D)) = P(X_A | \text{do}(X_B), X_C, \text{do}(X_D)) \quad \mu_{B \cup C}\text{-a.s.},$$

then it is sufficient to assume/check the following  $d$ -separation and absolute continuities:

$$A \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B | B \cup C \cup D, \quad \mu_{B \cup C} \ll P(X_B, X_C | \text{do}(X_D)) \ll \mu_{B \cup C},$$

$$\mu_C \ll P(X_C | \text{do}(X_B, X_D)) \ll \mu_C.$$

3. Insertion/deletion of action, for  $J \subseteq D$ : Assume that we want to establish the a.s.-equality:

$$P(X_A | \text{do}(X_B), X_C, \text{do}(X_D)) = P(X_A | X_C, \text{do}(X_D)) \quad \mu_C\text{-a.s.},$$

then it is sufficient to assume/check the following  $d$ -separation and absolute continuities:

$$A \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B | C \cup D, \quad \mu_C \ll P(X_C | \text{do}(X_B, X_D)) \ll \mu_C,$$

$$\mu_C \ll P(X_C | \text{do}(X_D)) \ll \mu_C.$$

4. Deletion of input: If

$$A \underset{G_{\text{do}(D)}}{\perp^d} J | C \cup D, \quad \mu_C \ll P(X_C | \text{do}(X_{D \cup J})) \ll \mu_C.$$

then there exists a Markov kernel  $P(X_A | X_C, \text{do}(X_D, \cancel{X_{J \setminus D}}))$  such that:

$$P(X_A | X_C, \text{do}(X_D, \cancel{X_{J \setminus D}})) = P(X_A | X_C, \text{do}(X_{D \cup J})) \quad \mu_C\text{-a.s.}$$

Note that the two-sided absolute continuities hold for  $\sigma$ -finite reference measures iff the indicated Markov kernel has a strictly positive Doob-Radon-Nikodym derivative w.r.t. the corresponding reference product measure by Corollary 2.3.20.

*Proof.* The proof follows directly from Theorem 5.1.2.

For the last rule (‘Deletion of input’), one can take the Markov kernel as

$$P(X_A|X_C, \text{do}(X_D, \cancel{X_{J \setminus D}})) := Q(X_A|X_C, X_D)$$

where  $Q(X_A|X_C, X_D)$  is defined in the proof of Proposition 5.1.7 point 3, for the special case  $B = I_B = \emptyset$ . The proof of Theorem 5.1.2 rule 3 then applies (as it doesn’t depend crucially on the assumption  $J \subseteq D$ , or  $B \neq \emptyset$ ), which shows the claim.  $\square$

**Remark 5.1.4.** *Note that in Corollary 5.1.3 (in contrast to Proposition 5.1.7 and Theorem 5.1.2) we cannot easily formulate the independence of variables  $X_{J \setminus D}$  in the presented way. If we accept the above then we can simplify the formulas (as done in Corollary 5.1.3) and assume that  $J \subseteq D$  and thus  $D \cup J = D$ , which makes then the d-separation requirements weaker (due to extra conditioning on  $J$ ). For the case where  $J$  does not fully lie in  $D$  one either needs to use Proposition 5.1.7 and Theorem 5.1.2 or the global Markov property, Theorem 4.2.1, directly.*

## Proofs - Do-Calculus

**Lemma 5.1.5.** *For pairwise disjoint  $B, C \subseteq V$  and  $D \subseteq V \cup J$  and a measurable subset  $N \subseteq \mathcal{X}_{B \cup C \cup D \cup J}$  the following statements are equivalent:*

1.  *$N$  is a  $P(X_B, X_C | \text{do}(X_{I_B}, X_{D \cup J}))$ -null set.*
2. *For every decomposition  $B = B_1 \dot{\cup} B_2$  the set  $N$  is a  $P(X_B, X_C | \text{do}(X_{B_2}, X_{D \cup J}))$ -null set.*
3. *For every decomposition  $B = B_1 \dot{\cup} B_2$  the set  $N$  is a  $P(X_{B_1}, X_C | \text{do}(X_{B_2}, X_{D \cup J}))$ -null set.*

*Proof.* Every value  $x_{I_B} = (x_{I_v})_{v \in B} \in \mathcal{X}_{I_B}$  defines a decomposition  $B = B_1 \dot{\cup} B_2$  via:

$$B_1 := \{v \in B \mid x_{I_v} = \star\}, \quad B_2 := \{v \in B \mid x_{I_v} \in \mathcal{X}_v\}.$$

So running through all values  $x_{I_B} \in \mathcal{X}_{I_B}$  is the same as running through all subsets  $B_2 \subseteq B$  and all values  $x_{B_2} \in \mathcal{X}_{B_2}$ , while putting  $x_{I_{B_1}} = \star$  for  $B_1 = B \setminus B_2$ . Furthermore, we have the following identities:

$$\begin{aligned} & P((X_B, X_C) \in N_{x_{D \cup J}} \mid \text{do}(X_{I_B} = x_{I_B}, X_{D \cup J} = x_{D \cup J})) \\ &= P((X_B, X_C) \in N_{x_{D \cup J}} \mid \text{do}(X_{I_{B_1}} = \star, X_{I_{B_2}} = x_{B_2}, X_{D \cup J} = x_{D \cup J})) \\ &= P((X_B, X_C) \in N_{x_{D \cup J}} \mid \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J} = x_{D \cup J})) \\ &= (P(X_{B_1}, X_C \mid \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J} = x_{D \cup J})) \otimes \delta(X_{B_2} \mid X_{B_2} = x_{B_2})) (N_{x_{D \cup J}}) \\ &= P((X_{B_1}, X_C) \in N_{(x_{B_2}, x_{D \cup J})} \mid \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J} = x_{D \cup J})). \end{aligned}$$

So the first line vanishes for all values  $x_{I_B} \in \mathcal{X}_{I_B}$  and  $x_{D \cup J} \in \mathcal{X}_{D \cup J}$  if and only if any other line vanishes for all subsets  $B_2 \subseteq B$  and all values  $x_{B_2} \in \mathcal{X}_{B_2}$  and  $x_{D \cup J} \in \mathcal{X}_{D \cup J}$ . This shows the claim.  $\square$

**Remark 5.1.6** (Null sets—again). *In the following we will often make statements like: “The Markov kernel  $K(X_A|X_B, X_C, X_D)$  is unique up to a measurable  $K(X_B|X_C, X_D)$ -null set in  $\mathcal{X}_{BUC}$ ”, (rather than in  $\mathcal{X}_{BUC \cup D}$ ). This means that the corresponding null set  $N$  can be considered constant in  $X_{D \setminus (BUC)}$ , or, more precisely, that  $N$  is of the form:*

$$N = M \times \mathcal{X}_{D \setminus (BUC)} \subseteq \mathcal{X}_{BUC \cup D},$$

with  $M \subseteq \mathcal{X}_{BUC}$ .

**Proposition 5.1.7** (Do-calculus—existence and uniqueness). *Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}G^+(v)})) \right)_{v \in V \cup U} \right).$$

Let  $A, B, C \subseteq V$  and  $D \subseteq J \cup V$  be such that  $A, B, C, D$  are pairwise disjoint. Then we have the following 3 rules relating marginal conditional to marginal interventional Markov kernels:

1. Insertion/deletion of observation: If we have:

$$A \underset{G_{\text{do}(D)}}{\overset{d}{\perp}} B | C \cup D,$$

then there exists a Markov kernel:

$$P(X_A | X_B, X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_C \times \mathcal{X}_D \dashrightarrow \mathcal{X}_A,$$

that is a version of:

$$P(X_A | X_{B_2}, X_C, \text{do}(X_{D \cup J})),$$

for every subset  $B_2 \subseteq B$  simultaneously. Note that this Markov kernel is only dependent on  $x_C$  and  $x_D$ , and constant in  $x_{J \setminus D}$ .

Such a Markov kernel is unique up to a measurable  $P(X_C | \text{do}(X_{D \cup J}))$ -null set in  $\mathcal{X}_{C \cup D}$ .

2. Action/observation exchange: If we have:

$$A \underset{G_{\text{do}(I_B, D)}}{\overset{d}{\perp}} I_B | B \cup C \cup D,$$

then there exists a Markov kernel:

$$P(X_A | \text{do}(X_B), X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_B \times \mathcal{X}_C \times \mathcal{X}_D \dashrightarrow \mathcal{X}_A,$$

that is a version of:

$$P(X_A | X_{B_1}, \text{do}(X_{B_2}), X_C, \text{do}(X_{D \cup J})),$$

for every decomposition:  $B = B_1 \dot{\cup} B_2$ , simultaneously.

Such a Markov kernel is unique up to a measurable  $P(X_B, X_C | \text{do}(X_{I_B}, X_{D \cup J}))$ -null set in  $N \subseteq \mathcal{X}_{B \cup C \cup D}$ , i.e.  $N$  is a  $P(X_{B_1}, X_C | \text{do}(X_{B_2}, X_{D \cup J}))$ -null set for every decomposition  $B = B_1 \dot{\cup} B_2$  simultaneously.

3. Insertion/deletion of action: If we have:

$$A \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B | C \cup D,$$

then there exists a Markov kernel:

$$P(X_A | \overline{\text{do}(X_B)}, X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_C \times \mathcal{X}_D \dashrightarrow \mathcal{X}_A,$$

that is a version of:

$$P(X_A | \text{do}(X_{B_2}), X_C, \text{do}(X_{D \cup J}))$$

for every subset  $B_2 \subseteq B$  simultaneously. Note that this Markov kernel is only dependent on  $x_C$  and  $x_D$ , and constant in  $x_{J \setminus D}$ .

Such a Markov kernel is unique up to a measurable  $P(X_C | \text{do}(X_{I_B}, X_{D \cup J}))$ -null set  $N \subseteq \mathcal{X}_{C \cup D}$ , i.e.  $N$  is a  $P(X_C | \text{do}(X_{B_2}, X_{D \cup J}))$ -null set for every subset  $B_2 \subseteq B$  simultaneously.

*Proof.* We make use of the global Markov property (GMP), theorem 4.2.1.

1.) The assumption:

$$A \underset{G_{\text{do}(D)}}{\perp^d} B | C \cup D,$$

implies the conditional independence by GMP 4.2.1:

$$X_A \underset{P(X_V | \text{do}(X_{D \cup J}))}{\perp\!\!\!\perp} X_B | X_C, X_D.$$

So we get the following factorization, where we can omit the deterministic variables from  $X_D$  on the left of the conditioning lines:

$$P(X_A, X_B, X_C | \text{do}(X_{D \cup J})) = Q(X_A | X_C, X_D) \otimes P(X_B, X_C | \text{do}(X_{D \cup J})),$$

for some Markov kernel  $Q(X_A | X_C, X_D)$ . Here  $Q(X_A | X_C, X_D)$  serves as a version of the conditional Markov kernel:

$$P(X_A | X_B, X_C, \text{do}(X_{D \cup J})).$$

If we marginalize out  $X_{B_1}$  for any decomposition  $B = B_1 \dot{\cup} B_2$  in the above factorization we also get:

$$P(X_A, X_C, X_{B_2} | \text{do}(X_{D \cup J})) = Q(X_A | X_C, X_D) \otimes P(X_C, X_{B_2} | \text{do}(X_{D \cup J})),$$

showing that  $Q(X_A | X_C, X_D)$  is also a version of:

$$P(X_A | X_{B_2}, X_C, \text{do}(X_{D \cup J})).$$

In particular, this holds for  $B_2 = \emptyset$ . This shows all the claimed properties for  $Q(X_A|X_C, X_D)$ .

Now consider another Markov kernel  $K(X_A|X_C, X_D)$  and the measurable sets:

$$\begin{aligned}\tilde{N} &:= \{x_{CUD} \in \mathcal{X}_{CUD} \mid Q(X_A|X_C = x_C, X_D = x_D) \neq K(X_A|X_C = x_C, X_D = x_D)\}, \\ N &:= \{x_{CUD \cup J} \in \mathcal{X}_{CUD \cup J} \mid Q(X_A|X_C = x_C, X_D = x_D) \neq K(X_A|X_C = x_C, X_D = x_D)\} \\ &= \tilde{N} \times \mathcal{X}_{J \setminus D}.\end{aligned}$$

If  $K(X_A|X_C, X_D)$  is a version of:

$$P(X_A|X_{B_2}, X_C, \text{do}(X_{D \cup J})),$$

for every subset  $B_2 \subseteq B$  simultaneously, then this holds, in particular, for  $B_2 = \emptyset$ . Since conditional Markov kernels are essentially unique, by Theorem 2.4.16, we have that  $N$  is a  $P(X_C|\text{do}(X_{D \cup J}))$ -null set.

2.) The assumption:

$$A \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B \mid B \cup C \cup D,$$

implies the conditional independence by GMP 4.2.1:

$$X_A \underset{P(X_V|\text{do}(X_{I_B}, X_{D \cup J}))}{\perp\!\!\!\perp} X_{I_B} \mid X_B, X_C, X_D.$$

So we have the following factorization:

$$P(X_A, X_B, X_C \mid \text{do}(X_{I_B}, X_{D \cup J})) = Q(X_A|X_B, X_C, X_D) \otimes P(X_B, X_C \mid \text{do}(X_{I_B}, X_{D \cup J})), \quad (21)$$

for some Markov kernel  $Q(X_A|X_B, X_C, X_D)$ , which serves as a version of the conditional Markov kernel:

$$P(X_A|X_B, X_C, \text{do}(X_{I_B}, X_{D \cup J})),$$

and which is independent of  $X_{I_B}$ .

We first claim that for a Markov kernel  $Q(X_A|X_B, X_C, X_D)$  the equation 21 is equivalent to the system of equations 22 indexed by subsets  $B_2 \subseteq B$  and with  $B_1 := B \setminus B_2$ :

$$P(X_A, X_{B_1}, X_C \mid \text{do}(X_{B_2}, X_{D \cup J})) = Q(X_A|X_B, X_C, X_D) \otimes P(X_{B_1}, X_C \mid \text{do}(X_{B_2}, X_{D \cup J})). \quad (22)$$

Indeed, we can look at the different input values for  $X_{I_B} = (X_{I_{B_1}}, X_{I_{B_2}})$  in equation 21. For  $B_1$  we put:  $X_{I_{B_1}} = \star = (\star)_{v \in B_1}$  and for  $B_2$  we take values:  $X_{I_{B_2}} = x_{B_2} \in \mathcal{X}_{B_2}$ . This implies:

$$\begin{aligned} &P(X_A, X_B, X_C \mid \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J})) \\ &= P(X_A, X_B, X_C \mid \text{do}(X_{I_{B_1}} = \star, X_{I_{B_2}} = x_{B_2}, X_{D \cup J})) \\ &= Q(X_A|X_B, X_C, X_D) \otimes P(X_B, X_C \mid \text{do}(X_{I_{B_1}} = \star, X_{I_{B_2}} = x_{B_2}, X_{D \cup J})), \\ &= Q(X_A|X_B, X_C, X_D) \otimes P(X_B, X_C \mid \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J})). \end{aligned}$$

So we get the equations:

$$P(X_A, X_B, X_C | \text{do}(X_{B_2}, X_{D \cup J})) = Q(X_A | X_B, X_C, X_D) \otimes P(X_B, X_C | \text{do}(X_{B_2}, X_{D \cup J})),$$

where we can further marginalize out the deterministic  $X_{B_2}$ :

$$P(X_A, X_{B_1}, X_C | \text{do}(X_{B_2}, X_{D \cup J})) = Q(X_A | X_B, X_C, X_D) \otimes P(X_{B_1}, X_C | \text{do}(X_{B_2}, X_{D \cup J})).$$

Note that we can also go back by multiplying with  $\delta(X_{B_2} | X_{B_2})$ . This shows the intermediate claim.

The equation 22 already implies that  $Q(X_A | X_B, X_C, X_D)$  is a version of the conditional Markov kernel:

$$P(X_A | X_{B_1}, X_C, \text{do}(X_{B_2}, X_{D \cup J})), \quad (23)$$

for every decomposition:  $B = B_1 \dot{\cup} B_2$  simultaneously.

Now consider another Markov kernel  $K(X_A | X_B, X_C, X_D)$  and the measurable sets:

$$\begin{aligned} \tilde{N} &:= \{x_{B \cup C \cup D} \in \mathcal{X}_{B \cup C \cup D} \mid Q(X_A | X_B = x_B, X_C = x_C, X_D = x_D) \\ &\quad \neq K(X_A | X_B = x_B, X_C = x_C, X_D = x_D)\}, \\ N &:= \tilde{N} \times \mathcal{X}_{J \setminus D}. \end{aligned}$$

If  $K(X_A | X_B, X_C, X_D)$  now is also a version of the conditional Markov kernel 23 for every decomposition  $B = B_1 \dot{\cup} B_2$  simultaneously, then  $N$  is a  $P(X_{B_1}, X_C | \text{do}(X_{B_2}, X_{D \cup J}))$ -null set for every decomposition  $B = B_1 \dot{\cup} B_2$ , because conditional Markov kernels are essentially unique, see Theorem 2.4.16. By Lemma 5.1.5 this statement is equivalent for  $N$  to be a  $P(X_B, X_C | \text{do}(X_{I_B}, X_{D \cup J}))$ -null set.

3.) The assumption:

$$A \underset{G_{\text{do}(I_B, D)}}{\overset{d}{\perp}} I_B \mid C \cup D,$$

implies the conditional independence by GMP 4.2.1:

$$X_A \underset{P(X_V | \text{do}(X_{I_B}, X_{D \cup J}))}{\perp\!\!\!\perp} X_{I_B} \mid X_C, X_D.$$

So we have the following factorization:

$$P(X_A, X_C | \text{do}(X_{I_B}, X_{D \cup J})) = Q(X_A | X_C, X_D) \otimes P(X_C | \text{do}(X_{I_B}, X_{D \cup J})), \quad (24)$$

for some Markov kernel  $Q(X_A | X_C, X_D)$ , which serves as a version of the conditional Markov kernel:

$$P(X_A | X_C, \text{do}(X_{I_B}, X_{D \cup J})),$$

and which is independent of  $X_{I_B}$ .

We can now look at the different input values for any decomposition:  $B = B_1 \dot{\cup} B_2$ . For this we put:  $X_{I_{B_1}} = \star = (\star)_{v \in B_1}$  and  $X_{I_{B_2}} = x_{B_2} \in \mathcal{X}_{B_2}$ . This implies:

$$\begin{aligned} & P(X_A, X_C | \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J})) \\ &= P(X_A, X_C | \text{do}(X_{I_{B_1}} = \star, X_{I_{B_2}} = x_{B_2}, X_{D \cup J})) \\ &= Q(X_A | X_C, X_D) \otimes P(X_C | \text{do}(X_{I_{B_1}} = \star, X_{I_{B_2}} = x_{B_2}, X_{D \cup J})), \\ &= Q(X_A | X_C, X_D) \otimes P(X_C | \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J})). \end{aligned}$$

So we get for every subset  $B_2 \subseteq B$ :

$$P(X_A, X_C | \text{do}(X_{B_2}, X_{D \cup J})) = Q(X_A | X_C, X_D) \otimes P(X_C | \text{do}(X_{B_2}, X_{D \cup J})), \quad (25)$$

which shows that  $Q(X_A | X_C, X_D)$  is a version of the conditional Markov kernel:

$$P(X_A | X_C, \text{do}(X_{B_2}, X_{D \cup J})), \quad (26)$$

for every subset  $B_2 \subseteq B$  simultaneously.

Now consider another Markov kernel  $K(X_A | X_C, X_D)$  and the measurable sets:

$$\begin{aligned} \tilde{N} &:= \{x_{C \cup D} \in \mathcal{X}_{C \cup D} | Q(X_A | X_C = x_C, X_D = x_D) \neq K(X_A | X_C = x_C, X_D = x_D)\}, \\ N_{B_2} &:= \mathcal{X}_{B_2} \times \tilde{N} \times \mathcal{X}_{J \setminus D}. \end{aligned}$$

Now assume that  $K(X_A | X_C, X_D)$  is a version of the conditional Markov kernel in (26) for every subset  $B_2 \subseteq B$  simultaneously. Then for every subset  $B_2 \subseteq B$  set  $N_{B_2}$  is a  $P(X_C | \text{do}(X_{B_2}, X_{D \cup J}))$ -null set, because of the essential uniqueness of conditional Markov kernels, see Theorem 2.4.16. More concretely, for  $x_{B_2} \in \mathcal{X}_{B_2}$  and  $x_{D \cup J} \in \mathcal{X}_{D \cup J}$  we get the equations:

$$\begin{aligned} 0 &= P(X_C \in (N_{B_2})_{(x_{B_2}, x_{D \cup J})} | \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J} = x_{D \cup J})) \\ &= P(X_C \in \tilde{N}_{x_D} | \text{do}(X_{B_2} = x_{B_2}, X_{D \cup J} = x_{D \cup J})) \\ &= P(X_C \in \tilde{N}_{x_D} | \text{do}(X_{I_{B_1}} = \star, X_{I_{B_2}} = x_{B_2}, X_{D \cup J} = x_{D \cup J})) \\ &= P(X_C \in \tilde{N}_{x_D} | \text{do}(X_{I_B} = (\star, x_{B_2}), X_{D \cup J} = x_{D \cup J})). \end{aligned}$$

Since we have this for all decompositions  $B = B_1 \dot{\cup} B_2$  and all values  $x_{B_2} \in \mathcal{X}_{B_2}$  we are running through all values  $x_{I_B} \in \mathcal{X}_{I_B}$ . This shows that  $\tilde{N}$  is a  $P(X_C | \text{do}(X_{I_B}, X_{D \cup J}))$ -null set in  $\mathcal{X}_{C \cup D}$ .  $\square$

**Theorem 5.1.2.** *[Almost-sure do-calculus—in detail] Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}G^+(v)})) \right)_{v \in V \cup U} \right).$$

Let  $A, B, C \subseteq V$  and  $D \subseteq J \cup V$  be such that  $A, B, C, D$  are pairwise disjoint.

Further assume that we have reference measures  $\mu_v$  on  $\mathcal{X}_v$  for every  $v \in V$  that are each equivalent to a probability measure (in terms of absolute continuity).<sup>23</sup> We then put  $\mu_F := \bigotimes_{v \in F} \mu_v$  for  $F \subseteq V$ .

<sup>23</sup>Recall the connection between absolute continuity and strictly positive densities in Corollary 2.3.20.

All  $\sigma$ -finite measures satisfy this assumption.



1. Insertion/deletion of observation: Assume:

$$A \underset{G_{\text{do}(D)}}{\perp}^d B \mid C \cup D.$$

For a fixed finite index set  $I$  consider subsets  $B^{(i)} \subseteq B$ , for  $i \in I$ , and pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A \mid X_{B^{(i)}}, X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{B \cup C \cup D \cup J} \rightarrow \mathcal{X}_{B^{(i)} \cup C \cup D \cup J} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_{B^{(i)}}, X_C \mid \text{do}(X_{D \cup J}))$ . Then there exists a measurable  $P(X_B, X_C \mid \text{do}(X_{D \cup J}))$ -null set  $N \subseteq \mathcal{X}_{B \cup C \cup D \cup J}$ , such that all those Markov kernels are equal on the complement  $N^c$ .

Note that if  $\mu_{B \cup C} \ll P(X_B, X_C \mid \text{do}(X_{D \cup J}))$  then  $N$  is also a  $\mu_{B \cup C}$ -null set, i.e. for every  $x_{D \cup J} \in \mathcal{X}_{D \cup J}$  we have:  $\mu_{B \cup C}(N_{x_{D \cup J}}) = 0$ .

If we also have the reverse  $P(X_B, X_C \mid \text{do}(X_{D \cup J})) \ll \mu_{B \cup C}$  then we can change the above conditional Markov kernels on a  $\mu_{B \cup C}$ -null set  $N$  while they remain versions of the corresponding conditional Markov kernel.<sup>24</sup>

2. Action/observation exchange: Assume:

$$A \underset{G_{\text{do}(I_B, D)}}{\perp}^d I_B \mid B \cup C \cup D.$$

For a fixed finite index set  $I$  consider decompositions  $B = B_1^{(i)} \dot{\cup} B_2^{(i)}$ , for  $i \in I$ , and pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A \mid X_{B_1^{(i)}}, \text{do}(X_{B_2^{(i)}}), X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{B \cup C \cup D \cup J} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_{B_1^{(i)}}, X_C \mid \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \otimes \mu_{B_2^{(i)}}$  and assume the following absolute continuities:

$$\mu_{B \cup C} \ll P(X_{B_1^{(i)}}, X_C \mid \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \otimes \mu_{B_2^{(i)}}$$

for all  $i \in I$ .<sup>25</sup> Then there exists a measurable  $\mu_{B \cup C}$ -null set  $N \subseteq \mathcal{X}_{B \cup C \cup D \cup J}$ , such that all those conditional Markov kernels are equal on the complement  $N^c$ .

<sup>24</sup>Note that the absolute continuities:  $\mu_{B \cup C} \ll P(X_B, X_C \mid \text{do}(X_{D \cup J})) \ll \mu_{B \cup C}$  hold if  $P(X_B, X_C \mid \text{do}(X_{D \cup J}))$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $\mu_{B \cup C}$ . Furthermore, the converse is also true for  $\sigma$ -finite reference measures  $\mu_{B \cup C}$  by Corollary 2.3.20.

<sup>25</sup>If you instead expected to pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A \mid X_{B_1^{(i)}}, \text{do}(X_{B_2^{(i)}}), X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{B \cup C \cup D \cup J} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_{B_1^{(i)}}, X_C \mid \text{do}(X_{B_2^{(i)}}), X_{D \cup J})$  and to assume the absolute continuities

$$\mu_{B_1^{(i)} \cup C} \ll P(X_{B_1^{(i)}}, X_C \mid \text{do}(X_{B_2^{(i)}}), X_{D \cup J}) \ll \mu_{B_1^{(i)} \cup C}$$

for all  $i \in I$ : that would lead to a similar, but slightly weaker statement.

If we also assume the reverse absolute continuities for all  $i \in I$ :

$$P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}, X_{D \cup J})) \otimes \mu_{B_2^{(i)}} \ll \mu_{BUC},$$

then all those conditional Markov kernels are versions of each other.<sup>26</sup>

3. Insertion/deletion of action: Assume:

$$A \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B | C \cup D.$$

For a fixed finite index set  $I$  consider subsets  $B^{(i)} \subseteq B$ , for  $i \in I$ , and pick for each  $i \in I$  an arbitrary version of a conditional Markov kernel:

$$P(X_A | \text{do}(X_{B^{(i)}}), X_C, \text{do}(X_{D \cup J})) : \mathcal{X}_{BUCUDUJ} \rightarrow \mathcal{X}_{B^{(i)}UCUDUJ} \rightarrow \mathcal{P}(\mathcal{X}_A),$$

of  $P(X_A, X_C | \text{do}(X_{B^{(i)}}, X_{D \cup J}))$  and assume the following absolute continuities:

$$\mu_C \ll P(X_C | \text{do}(X_{B^{(i)}}, X_{D \cup J}))$$

for all  $i \in I$ . Then there exists a measurable  $\mu_C$ -null set  $N \subseteq \mathcal{X}_{BUCUDUJ}$ , such that all those conditional Markov kernels are equal on the complement  $N^c$ .

If we also assume the reverse absolute continuities for all  $i \in I$ :

$$P(X_C | \text{do}(X_{B^{(i)}}, X_{D \cup J})) \ll \mu_C,$$

then all those conditional Markov kernels are versions of each other.<sup>27</sup>

*Proof.* W.l.o.g. we can assume all  $\mu_v$  to be probability measures.

1.) Let  $Q(X_A | X_C, X_D)$  be the Markov kernel from Proposition 5.1.7 point 1. Recall that conditional Markov kernels are essentially unique by Theorem 2.4.16. This shows that the set:

$$\begin{aligned} \tilde{N}^{(i)} := \{ & x_{B^{(i)}UCUDUJ} \in \mathcal{X}_{B^{(i)}UCUDUJ} \mid Q(X_A | X_C = x_C, X_D = x_D) \\ & \neq P(X_A | X_{B^{(i)}} = x_{B^{(i)}}, X_C = x_C, \text{do}(X_{D \cup J} = x_{D \cup J})) \}, \end{aligned}$$

is a (measurable)  $P(X_{B^{(i)}}, X_C | \text{do}(X_{D \cup J}))$ -null set. So the lifted set:

$$\begin{aligned} N^{(i)} := \{ & x_{BUCUDUJ} \in \mathcal{X}_{BUCUDUJ} \mid Q(X_A | X_C = x_C, X_D = x_D) \\ & \neq P(X_A | X_{B^{(i)}} = x_{B^{(i)}}, X_C = x_C, \text{do}(X_{D \cup J} = x_{D \cup J})) \}, \end{aligned}$$

<sup>26</sup>Note that the absolute continuities:  $\mu_{BUC} \ll P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}, X_{D \cup J})) \otimes \mu_{B_2^{(i)}} \ll \mu_{BUC}$  hold if the absolute continuities:  $\mu_{B_1^{(i)}UC} \ll P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}, X_{D \cup J})) \ll \mu_{B_1^{(i)}UC}$  hold, which hold if  $P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}, X_{D \cup J}))$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $\mu_{B_1^{(i)}UC}$ . Furthermore, the converse is also true for  $\sigma$ -finite reference measures  $\mu_{B_1^{(i)}UC}$  by Corollary 2.3.20.

<sup>27</sup>Note that absolute continuities:  $\mu_C \ll P(X_C | \text{do}(X_{B^{(i)}}, X_{D \cup J})) \ll \mu_C$  hold if  $P(X_C | \text{do}(X_{B^{(i)}}, X_{D \cup J}))$  has a strictly positive Doob-Radon-Nikodym derivative w.r.t.  $\mu_C$ . Furthermore, the converse is also true for  $\sigma$ -finite reference measures  $\mu_C$  by Corollary 2.3.20.

is then a (measurable)  $P(X_B, X_C | \text{do}(X_{D \cup J}))$ -null set. Then also the finite union:

$$N := \bigcup_{i \in I} N^{(i)} \subseteq \mathcal{X}_{B \cup C \cup D \cup J},$$

is a (measurable)  $P(X_B, X_C | \text{do}(X_{D \cup J}))$ -null set as well. Note that on the complement  $N^c$  all Markov kernels agree with  $Q(X_A | X_C, X_D)$  and are thus all equal on  $N^c$ .

2.) Consider the Markov kernel  $Q(X_A | X_B, X_C, X_D)$  from Proposition 5.1.7 point 2 and for  $i \in I$  the measurable set:

$$N^{(i)} := \left\{ x_{B \cup C \cup D \cup J} \in \mathcal{X}_{B \cup C \cup D \cup J} \mid \begin{aligned} &Q(X_A | X_B = x_B, X_C = x_C, X_D = x_D) \\ &\neq P(X_A | X_{B_1^{(i)}} = x_{B_1^{(i)}}, \text{do}(X_{B_2^{(i)}} = x_{B_2^{(i)}}), X_C = x_C, \text{do}(X_{D \cup J} = x_{D \cup J})) \end{aligned} \right\}.$$

Again, by the essential uniqueness of conditional Markov kernels, Theorem 2.4.16, the set  $N^{(i)}$  is a  $P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}, X_{D \cup J})) \otimes \mu_{B_2^{(i)}}$ -null set. The absolute continuity:

$$\mu_{B \cup C} \ll P(X_{B_1^{(i)}}, X_C | \text{do}(X_{B_2^{(i)}}, X_{D \cup J})) \otimes \mu_{B_2^{(i)}},$$

then renders  $N^{(i)}$  a  $\mu_{B \cup C}$ -null set. This shows that the finite union:

$$N := \bigcup_{i \in I} N^{(i)},$$

is a  $\mu_{B \cup C}$ -null set as well. Again, note that on the complement  $N^c$  all Markov kernels agree with  $Q(X_A | X_B, X_C, X_D)$  and are thus all equal on  $N^c$ .

3.) Consider the Markov kernel  $Q(X_A | X_C, X_D)$  from Proposition 5.1.7 point 3 and for  $i \in I$  the measurable set:

$$\tilde{N}^{(i)} := \left\{ x_{B^{(i)} \cup C \cup D \cup J} \in \mathcal{X}_{B^{(i)} \cup C \cup D \cup J} \mid \begin{aligned} &Q(X_A | X_C = x_C, X_D = x_D) \\ &\neq P(X_A | \text{do}(X_{B^{(i)}} = x_{B^{(i)}}), X_C = x_C, \text{do}(X_{D \cup J} = x_{D \cup J})) \end{aligned} \right\}.$$

Again, by the essential uniqueness of conditional Markov kernels, Theorem 2.4.16, the set  $\tilde{N}^{(i)}$  is a  $P(X_C | \text{do}(X_{B^{(i)}}, X_{D \cup J}))$ -null set. By the absolute continuity  $\mu_C \ll P(X_C | \text{do}(X_{B^{(i)}}, X_{D \cup J}))$  we get that  $\tilde{N}^{(i)}$  is a  $\mu_C$ -null set. This shows that the measurable set:

$$N^{(i)} := \left\{ x_{B \cup C \cup D \cup J} \in \mathcal{X}_{B \cup C \cup D \cup J} \mid \begin{aligned} &Q(X_A | X_C = x_C, X_D = x_D) \\ &\neq P(X_A | \text{do}(X_{B^{(i)}} = x_{B^{(i)}}), X_C = x_C, \text{do}(X_{D \cup J} = x_{D \cup J})) \end{aligned} \right\}.$$

is a  $\mu_C$ -null set as well. Then the finite union:

$$N := \bigcup_{i \in I} N^{(i)},$$

is also a  $\mu_C$ -null set. Again, note that on the complement  $N^c$  all Markov kernels agree with  $Q(X_A | X_C, X_D)$  and are thus all equal on  $N^c$ .  $\square$

## 5.2. Adjustment Criteria and Formulae

**Motivation 5.2.1.** Consider an L-CBN:

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}^{G^+}(v)})) \right)_{v \in V \cup U} \right).$$

For simplicity assume that there are no input variables, i.e.  $J = \emptyset$ . Then the joint distribution is “do-free” and given as:

$$P(X_V, X_U) = \bigotimes_{v \in U \cup V} P_v \left( X_v | \text{do}(X_{\text{Pa}^{G^+}(v)}) \right),$$

with observational distribution as its marginal:  $P(X_V)$ .

We also have all the interventional distributions for  $W \subseteq V$ :

$$P(X_{V \setminus W}, X_U | \text{do}(X_W)) = \bigotimes_{v \in U \cup V \setminus W} P_v \left( X_v | \text{do}(X_{\text{Pa}^{G^+}(v)}) \right),$$

with marginals:  $P(X_{V \setminus W} | \text{do}(X_W))$ .

If we wanted to learn the distribution  $P(X_V)$  we could do an observational study and apply the usual statistical or machine learning techniques. If, in contrast, we wanted to learn interventional distributions:  $P(X_{V \setminus W} | \text{do}(X_W))$  from data (e.g. whether vaccination makes people immune to a disease), we typically would need to perform an interventional study where we intervene on the variables  $X_W$  and set them to different values. This usually requires expensive, time-consuming randomized control trials with an own group for each possible value of  $X_W$ .

If we assume that we know the causal graph  $G^+$  or  $G$  we could try to leverage the rules of do-calculus in a clever way and might be able to go from expressions involving  $\text{do}(W)$  to expressions only involving  $\text{do}(D)$  for a (much) smaller subset  $D \subseteq W$ , ideally  $D = \emptyset$ . Practically this would mean that we would need a much smaller randomized control trial and save time and resources.

For example, if we have the graph only involving the edge:  $v_1 \rightarrow v_2$  we have that:

$$P(X_2 | \text{do}(X_1)) = P(X_2 | X_1),$$

which can be estimated using observational data only, e.g. via supervised learning.

So the question of identifiability is now: Assuming that the causal graph is known, under which circumstances is a causal effect  $P(X_A | \text{do}(X_B))$  already determined by the observational distribution  $P(X_V)$ ? When can causal effects be identified via distributions that have less interventions in them?

**Notation 5.2.2.** Consider an L-CBN:

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}^{G^+}(v)})) \right)_{v \in V \cup U} \right).$$

We are interested in estimating the conditional causal effect:

$$P(X_A | X_C, \text{do}(X_B, X_D)),$$

but we only have data from:

$$P(X_V|X_C, \text{do}(X_D)).$$

The following index sets will have the following roles:

1.  $A$ : the outcome variables of interest.
2.  $B$ : the treatment or intervention variables.
3.  $C$ : general conditional (context) variables under which the data was collected.
4.  $D$ : general interventional (context) variables that were set by the experimenter,  $J \subseteq D$ .
5.  $F_0$ : core adjustment variables, i.e. features that were measured.
6.  $F_1$ : additional measured adjustment variables.
7.  $F = F_0 \cup F_1$ .
8.  $H$ : additional unobserved variables.

**Theorem 5.2.3** (General adjustment formula). *Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}G^+(v)})) \right)_{v \in V \cup U} \right).$$

Assume that all the following conditions hold in the graphs  $G_{\text{do}(I_B, D)}^+$ :

$$(F_0 \cup H) \underset{G_{\text{do}(I_B, D)}^+}{\perp^d} I_B | (C \cup D), \quad (27)$$

$$A \underset{G_{\text{do}(I_B, D)}^+}{\perp^d} (F_1 \cup I_B) | (B \cup F_0 \cup H \cup C \cup D), \quad (28)$$

$$H \underset{G_{\text{do}(I_B, D)}^+}{\perp^d} B | (F \cup C \cup I_B \cup D). \quad (29)$$

Further assume that we have reference measures  $\mu_v$  on  $\mathcal{X}_v$ ,  $v \in V \cup H$ , such that:

$$\begin{aligned} \mu_{B \cup C \cup F \cup H} &\ll P(X_B, X_C, X_F, X_H | \text{do}(X_D)) \ll \mu_{B \cup C \cup F \cup H}, \\ \mu_{C \cup F \cup H} &\ll P(X_C, X_F, X_H | \text{do}(X_B, X_D)) \ll \mu_{C \cup F \cup H}. \end{aligned}$$

Then we have the adjustment formula:

$$P(X_A | X_C, \text{do}(X_B, X_D)) = P(X_A | X_B, X_C, X_F, \text{do}(X_D)) \circ P(X_F | X_C, \text{do}(X_D)) \quad \mu_{B \cup C}\text{-a.s.}$$

*Proof.* With help of Corollary 5.1.3 (2nd rule) we can establish the a.s.-equality:

$$P(X_A|X_{F_0}, X_H, X_C, X_B, \text{do}(X_D)) = P(X_A|X_{F_0}, X_H, X_C, \text{do}(X_B, X_D)) \quad \mu_{F_0 \cup H \cup C \cup B}\text{-a.s.}, \quad (30)$$

using the assumptions (implied by eq. 28):

$$\begin{aligned} A & \stackrel{d}{\perp}_{G_{\text{do}(I_B, D)}^+} I_B | (B \cup F_0 \cup H \cup C \cup D), \\ \mu_{F_0 \cup H \cup C \cup B} & \ll P(X_{F_0}, X_H, X_C, X_B | \text{do}(X_D)) \ll \mu_{F_0 \cup H \cup C \cup B}, \\ \mu_{F_0 \cup H \cup C} & \ll P(X_{F_0}, X_H, X_C | \text{do}(X_B, X_D)) \ll \mu_{F_0 \cup H \cup C}. \end{aligned}$$

With help of Corollary 5.1.3 (3rd rule) we can establish the a.s.-equality:

$$P(X_{F_0}, X_H | X_C, \text{do}(X_B, X_D)) = P(X_{F_0}, X_H | X_C, \text{do}(X_D)) \quad \mu_C\text{-a.s.}, \quad (31)$$

using the assumptions (implied by eq. 27):

$$\begin{aligned} (F_0 \cup H) & \stackrel{d}{\perp}_{G_{\text{do}(I_B, D)}^+} I_B | (C \cup D), \\ \mu_C & \ll P(X_C | \text{do}(X_B, X_D)) \ll \mu_C, \\ \mu_C & \ll P(X_C | \text{do}(X_D)) \ll \mu_C. \end{aligned}$$

With help of Corollary 5.1.3 (1st rule) we can establish the a.s.-equality:

$$P(X_A|X_{F_0}, X_H, X_C, X_B, \text{do}(X_D)) = P(X_A|X_{F_0}, X_{F_1}, X_H, X_C, X_B, \text{do}(X_D)) \quad \mu_{F_0 \cup F_1 \cup H \cup C \cup B}\text{-a.s.}, \quad (32)$$

using the assumptions (implied by eq. 28):

$$\begin{aligned} A & \stackrel{d}{\perp}_{G_{\text{do}(D)}^+} F_1 | (B \cup F_0 \cup H \cup C \cup D), \\ \mu_{F_0 \cup F_1 \cup H \cup C \cup B} & \ll P(X_{F_0}, X_{F_1}, X_H, X_C, X_B | \text{do}(X_D)) \ll \mu_{F_0 \cup F_1 \cup H \cup C \cup B}. \end{aligned}$$

With help of Corollary 5.1.3 (1st rule) we can establish the a.s.-equality:

$$P(X_H|X_F, X_C, \text{do}(X_D)) = P(X_H|X_F, X_C, X_B, \text{do}(X_D)) \quad \mu_{F \cup C \cup B}\text{-a.s.}, \quad (33)$$

using the assumptions (implied by eq. 29):

$$\begin{aligned} H & \stackrel{d}{\perp}_{G_{\text{do}(D)}^+} B | (F \cup C \cup D), \\ \mu_{F \cup C \cup B} & \ll P(X_F, X_C, X_B | \text{do}(X_D)) \ll \mu_{F \cup C \cup B}. \end{aligned}$$

These a.s.-equations together with the chain rule gives us the following  $\mu_{BUC}$ -a.s.-equation:

$$\begin{aligned}
& P(X_A|X_C, \text{do}(X_B, X_D)) \\
&= P(X_A|X_{F_0}, X_H, X_C, \text{do}(X_B, X_D)) \circ P(X_{F_0}, X_H|X_C, \text{do}(X_B, X_D)) \\
&\stackrel{30}{=} P(X_A|X_{F_0}, X_H, X_C, X_B, \text{do}(X_D)) \circ P(X_{F_0}, X_H|X_C, \text{do}(X_B, X_D)) \\
&\stackrel{31}{=} P(X_A|X_{F_0}, X_H, X_C, X_B, \text{do}(X_D)) \circ P(X_{F_0}, X_H|X_C, \text{do}(X_D)) \\
&= P(X_A|X_{F_0}, X_H, X_C, X_B, \text{do}(X_D)) \circ P(X_{F_0}, X_{F_1}, X_H|X_C, \text{do}(X_D)) \\
&\stackrel{32}{=} P(X_A|X_{F_0}, X_{F_1}, X_H, X_C, X_B, \text{do}(X_D)) \circ P(X_{F_0}, X_{F_1}, X_H|X_C, \text{do}(X_D)) \\
&= P(X_A|X_F, X_H, X_C, X_B, \text{do}(X_D)) \circ P(X_F, X_H|X_C, \text{do}(X_D)) \\
&= P(X_A|X_F, X_H, X_C, X_B, \text{do}(X_D)) \circ (P(X_H|X_F, X_C, \text{do}(X_D)) \otimes P(X_F|X_C, \text{do}(X_D))) \\
&\stackrel{33}{=} P(X_A|X_F, X_H, X_C, X_B, \text{do}(X_D)) \circ (P(X_H|X_F, X_C, X_B, \text{do}(X_D)) \otimes P(X_F|X_C, \text{do}(X_D))) \\
&= P(X_A|X_F, X_C, X_B, \text{do}(X_D)) \circ P(X_F|X_C, \text{do}(X_D)).
\end{aligned}$$

Note that the disintegration:

$$P(X_F, X_H|X_C, \text{do}(X_D)) = P(X_H|X_F, X_C, \text{do}(X_D)) \otimes P(X_F|X_C, \text{do}(X_D))$$

holds (only)  $P(X_C|\text{do}(X_D))$ -a.s., as for the conditional  $P(X_H|X_F, X_C, \text{do}(X_D))$  we have the ambiguity if it is considered a conditional of  $P(X_F, X_H|X_C, \text{do}(X_D))$ , for which then we have “sure” equality, or, if it is considered a conditional of  $P(X_F, X_H, X_C, \text{do}(X_D))$ , for which then we have only the above “almost-sure” equality. Further note, that by the assumption on the reference measures, the above equality then also holds  $\mu_C$ -a.s.  $\square$

**Corollary 5.2.4** (Conditional interventional backdoor covariate adjustment formula). *Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}G^+(v)})) \right)_{v \in V \cup U} \right).$$

Assume that the conditional interventional backdoor criterion in the graphs  $G_{\text{do}(I_B, D)}$  holds:

1.  $F \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B | (C \cup D)$ , and:
2.  $A \underset{G_{\text{do}(I_B, D)}}{\perp^d} I_B | (B \cup F \cup C \cup D)$ .

Further assume the following absolute continuities:

$$\begin{aligned}
\mu_{BUC \cup F} &\ll P(X_B, X_C, X_F | \text{do}(X_D)) \ll \mu_{BUC \cup F}, \\
\mu_{C \cup F} &\ll P(X_C, X_F | \text{do}(X_B, X_D)) \ll \mu_{C \cup F}.
\end{aligned}$$

Then we have the adjustment formula:

$$P(X_A|X_C, \text{do}(X_B, X_D)) = P(X_A|X_B, X_C, X_F, \text{do}(X_D)) \circ P(X_F|X_C, \text{do}(X_D)) \quad \mu_{BUC}\text{-a.s.}$$

*Proof.* It follows by the same arguments as in Theorem 5.2.3 with  $F_1 = H = \emptyset$ .  $\square$

Without the conditioning set, i.e.  $C = \emptyset$ , and direct careful analysis we get a version with slightly weaker positivity assumptions:

**Corollary 5.2.5** (Interventional backdoor covariate adjustment formula). *Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}G^+(v)})) \right)_{v \in V \cup U} \right).$$

Assume that the interventional backdoor criterion in the graphs  $G_{\text{do}(I_B, D)}$  holds:

1.  $F \stackrel{d}{\perp}_{G_{\text{do}(I_B, D)}} I_B | D$ , and:
2.  $A \stackrel{d}{\perp}_{G_{\text{do}(I_B, D)}} I_B | (B \cup F \cup D)$ .

Further assume the following absolute continuity:

$$P(X_F | \text{do}(X_D)) \otimes P(X_B | \text{do}(X_D)) \ll P(X_F, X_B | \text{do}(X_D)).$$

Then we have the adjustment formulas:

$$\begin{aligned} P(X_A, X_F | \text{do}(X_B, X_D)) &= P(X_A | X_F, X_B, \text{do}(X_D)) \otimes P(X_F | \text{do}(X_D)) \quad P(X_B | \text{do}(X_D))\text{-a.s.}, \\ P(X_A | \text{do}(X_B, X_D)) &= P(X_A | X_F, X_B, \text{do}(X_D)) \circ P(X_F | \text{do}(X_D)) \quad P(X_B | \text{do}(X_D))\text{-a.s.} \end{aligned}$$

*Proof.* First, note that “do( $X_D$ )” appears in every Markov kernel in the above formulas. So for readability, we will drop it in the following everywhere.

By the first d-separation assumption we see by Proposition 5.1.7 (rule 3) that we have the “sure” equality:  $P(X_F | \text{do}(X_B)) = P(X_F)$ . By the second d-separation assumption we see by Proposition 5.1.7 (rule 2) that we have a Markov kernel  $P(X_A | X_F, \text{do}(X_B))$  that is also a version of  $P(X_A | X_F, X_B)$ . So any version of  $P(X_A | X_F, X_B)$  can be changed on a  $P(X_F, X_B)$ -null set  $N \subseteq \mathcal{X}_B \times \mathcal{X}_F$  to get  $P(X_A | X_F, \text{do}(X_B))$ . The absolute continuity assumption implies that  $N$  is also a  $P(X_F) \otimes P(X_B)$ -null set. This implies that we have the equations of Markov kernels:

$$\begin{aligned} P(X_A | X_F, X_B) \otimes P(X_F) \otimes P(X_B) &= P(X_A | X_F, \text{do}(X_B)) \otimes P(X_F) \otimes P(X_B) \\ &= P(X_A | X_F, \text{do}(X_B)) \otimes P(X_F | \text{do}(X_B)) \otimes P(X_B) \\ &= P(X_A, X_F | \text{do}(X_B)) \otimes P(X_B). \end{aligned}$$

By the essential uniqueness of conditional Markov kernels we get that:

$$P(X_A, X_F | \text{do}(X_B)) = P(X_A | X_F, X_B) \otimes P(X_F) \quad P(X_B)\text{-a.s.}$$

Marginalizing out  $X_F$  on both sides gives us the remaining claim.  $\square$

We can now further specialize to the case with  $C = D = J = \emptyset$  and immediately get:



**Corollary 5.2.6** (Backdoor covariate adjustment). *Assume that the backdoor criterion holds:*

1.  $F \perp_{G_{\text{do}(I_B)}}^d I_B$ , and:
2.  $A \perp_{G_{\text{do}(I_B)}}^d I_B | (B \cup F)$ .

Further assume the following absolute continuity:

$$P(X_F) \otimes P(X_B) \ll P(X_F, X_B).$$

Then we have the adjustment formulas:

$$\begin{aligned} P(X_A, X_F | \text{do}(X_B)) &= P(X_A | X_F, X_B) \otimes P(X_F) && P(X_B)\text{-a.s.}, \\ P(X_A | \text{do}(X_B)) &= P(X_A | X_F, X_B) \circ P(X_F) && P(X_B)\text{-a.s.} \end{aligned}$$

**Remark 5.2.7.** *An example how the adjustment formula may fail if the strict positivity assumptions are not met is provided in Example 5.3.29.*

### 5.3. The ID-Algorithm

Consider an L-CBN:

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}G^+(v)})) \right)_{v \in V \cup U} \right),$$

with observable CADMG  $G$ . For subsets  $A, B, C \subseteq V$  we want to infer the conditional interventional distribution  $P(X_A | X_B, \text{do}(X_{J \cup C}))$  in terms of (repeated products and) conditional marginals of the observable Markov kernel  $P(X_V | \text{do}(X_J))$  and knowledge of  $G$ . In this subsection we will restrict ourselves to the case  $P(X_A | \text{do}(X_{J \cup C}))$  with no conditioning and present the *ID-algorithm*, which can tell if this is possible or not (in precise terms), and if so, provides us with a formula to do so. We will start this subsection with a series of necessary definitions, notations and lemmata. The main references are [Pea09, GP95, Tia02, TP02, Tia04, SP06b, HV06, HV08, RERS23, FM20].

#### 5.3.1. Core Definitions and Notations

**Definition 5.3.1** (Identifiability of interventional distributions/Markov kernels). *Let  $G = (J, V, E, L)$  be a CADMG and  $B \subseteq V$  and  $C \subseteq J \cup V$  disjoint subsets. We say that the interventional distribution/Markov kernel of  $C$  onto  $B$  is identifiable from  $G$ , or, more in generic symbols, that  $P(X_B | \text{do}(X_{J \cup C}))$  is identifiable from  $P(X_V | \text{do}(X_J))$  (and  $G$ ), if for every two L-CBNs  $M_1$  and  $M_2$  with the same:*

1. *observable CADMG  $G_1 = G_2 = G$ , and:*
2. *underlying spaces  $\mathcal{X}_{1,v} = \mathcal{X}_{2,v} =: \mathcal{X}_v$  for  $v \in J \cup V$ , and:*

3. *observable Markov kernels*  $P_1(X_V | \text{do}(X_J)) = P_2(X_V | \text{do}(X_J))$ ,

we also have the equality of the interventional Markov kernels:

$$P_1(X_B | \text{do}(X_{J \cup C})) = P_2(X_B | \text{do}(X_{J \cup C})).$$

Sometimes we further restrict the class of CBNs to define/achieve identifiability, e.g. by adding “for linear Gaussian CBNs” or “for discrete CBNs with strictly positive mass functions”, etc., and then only require  $M_1$  and  $M_2$  to come from such classes.

In the following we introduce a somewhat more vague, but constructive, notion of identifiability, which we coin *trackability* that allows us to follow certain marginalization, conditioning and multiplication steps to arrive at the wanted interventional Markov kernel.

**Definition 5.3.2** (Trackability (up to specifications)). *Let  $G = (J, V, E, L)$  be a CADMG and  $B \subseteq V$  and  $C \subseteq J \cup V$  disjoint subsets.*

1. *We say that the interventional distribution/Markov kernel of  $C$  onto  $B$  is trackable from  $G$ , or simply that  $P(X_B | \text{do}(X_{J \cup C}))$  is trackable from  $P(X_V | \text{do}(X_J))$  (and  $G$ ), if there exists a finite sequence of operations, only involving marginalization, conditioning and multiplication of Markov kernels, applied to previously determined Markov kernels, starting from  $P(X_V | \text{do}(X_J))$ , with predetermined target sets  $T_n \subseteq G$ , indicating on which variable the operation is applied to, such that for every  $L$ -CBNs  $M$  with observable CADMG  $G$  we can compute  $P(X_B | \text{do}(X_{J \cup C}))$  from  $P(X_V | \text{do}(X_J))$  when we follow the above sequence of operations (and this should work no matter which version of conditional Markov kernels were used).*
2. *We say that  $P(X_B | \text{do}(X_{J \cup C}))$  is trackable up to specifications “xyz” from  $P(X_V | \text{do}(X_J))$  (and  $G$ ), if the same as above holds true, but whenever we condition, which leads to a Markov kernel only up to some null sets, we use the specifications “xyz” to pick a certain version of conditional Markov kernel at each step such that following the pre-specified operations leads to  $P(X_B | \text{do}(X_{J \cup C}))$ .*
3. *We say that  $P(X_B | \text{do}(X_{J \cup C}))$  is trackable up to oracle choices from  $P(X_V | \text{do}(X_J))$  (and  $G$ ) if there exists a conditional Markov kernel at each conditioning step (“chosen by an oracle that knows  $M$ ”) such that following these operations leads to  $P(X_B | \text{do}(X_{J \cup C}))$ .*
4. *Again, we sometimes further restrict the class of CBNs to define/achieve trackability (up to oracle choices), e.g. by adding “for linear Gaussian CBNs” or “for discrete CBNs with strictly positive mass functions”, etc., and then only allow  $M$  to come from such classes.*

**Example 5.3.3.** *To illustrate how such a series of operations could look like consider the CADMG  $G$  from Figure 11 with  $V = \{v_1, v_2, v_3\}$ . We assume that the observational distribution  $P(X_1, X_2, X_3)$  is given. A list of operations could look like:*

1. Condition  $P(X_1, X_2, X_3)$  on  $(X_1, X_2)$  and get a version  $P(X_3|X_1, X_2)$ .  
Further specification could be (if possible): “take a continuous version” or “take a version that is only dependent on variables  $X_1$ ” or “take a strictly positive version”.
2. Marginalize out  $(X_2, X_3)$  from  $P(X_1, X_2, X_3)$  and get  $P(X_1)$ .
3. Take the product of the previous two Markov kernels:  $P(X_3|X_1, X_2) \otimes P(X_1)$ .
4. Marginalize out  $X_1$  from the last Markov kernel and get:  $P(X_3|X_1, X_2) \circ P(X_1)$ .

**Lemma 5.3.4.** Let  $G = (J, V, E, L)$  be a CADMG.

1. If  $P(X_B | \text{do}(X_{J \cup C}))$  is trackable from  $P(X_V | \text{do}(X_J))$  then it is also identifiable.
2. If  $P(X_B | \text{do}(X_{J \cup C}))$  is trackable up to oracle choices from  $P(X_V | \text{do}(X_J))$  then it is also trackable (and thus identifiable) for discrete CBNs  $M$  with strictly positive mass functions:  $p(x_V | \text{do}(x_J)) > 0$  for all  $x_V, x_J$ .

*Proof.* The first point is clear as the sequence of operations always ends in the same result. For discrete CBNs  $M$  with strictly positive mass functions conditional Markov kernels are unambiguous, thus a sequence of marginalization, conditioning and products always leads to the same result. Note that marginals, conditionals and products of strictly positive mass functions also are strictly positive mass functions.  $\square$

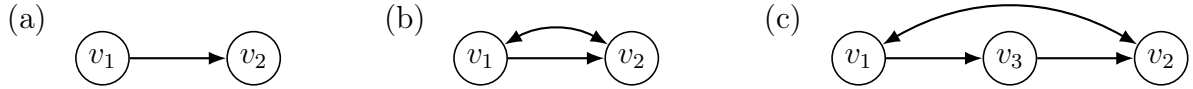


Figure 10: (a) A DAG with two nodes. (b) An ADMG with two nodes. (c) An ADMG with three nodes. The interventional distribution  $P(X_2 | \text{do}(X_1))$  is trackable up to oracle choices from  $P(X_1, X_2, X_3)$  in (a) and (c), but not in (b).

**Example 5.3.5.** Consider the DAG  $G = (V, E)$  from Figure 10 (a) with  $V = \{v_1, v_2\}$  and  $E = \{v_1 \rightarrow v_2\}$ . Let  $\mathcal{X}_1 := \{a, b, c\}$  and  $\mathcal{X}_2 := \{0, 1\}$ . Define the following Markov kernels  $P_1(X_1) := P_1(X_2) := P(X_1)$  via:

$$P(X_1 = a) := \frac{1}{2}, \quad P(X_1 = b) := \frac{1}{2}, \quad P(X_1 = c) := 0,$$

and further:

$$\begin{aligned} P_1(X_2 | \text{do}(X_1 = a)) &= \text{Bern}(1/4), & P_2(X_2 | \text{do}(X_1 = a)) &= \text{Bern}(1/4), \\ P_1(X_2 | \text{do}(X_1 = b)) &= \text{Bern}(3/4), & P_2(X_2 | \text{do}(X_1 = b)) &= \text{Bern}(3/4), \\ P_1(X_2 | \text{do}(X_1 = c)) &= \text{Bern}(1/8), & P_2(X_2 | \text{do}(X_1 = c)) &= \text{Bern}(7/8). \end{aligned}$$

Then we have two CBNs with observable DAG  $G$ , the same underlying spaces and the same observable Markov kernel:

$$P(X_1, X_2) := P_1(X_2 | \text{do}(X_1)) \otimes P_1(X_1) = P_2(X_2 | \text{do}(X_1)) \otimes P_2(X_1),$$

given by:

$$\begin{aligned} M_1 &:= (G, (P_1(X_1), P_1(X_2 | \text{do}(X_1)))) , \\ M_2 &:= (G, (P_2(X_1), P_2(X_2 | \text{do}(X_1)))) . \end{aligned}$$

Furthermore, consider the Markov kernel  $P(X_2 | X_1)$  given by:

$$\begin{aligned} P(X_2 | X_1 = a) &= \text{Bern}(1/4), \\ P(X_2 | X_1 = b) &= \text{Bern}(3/4), \\ P(X_2 | X_1 = c) &= \text{Bern}(5/8). \end{aligned}$$

We thus have three different versions of the conditional Markov kernels of  $P(X_1, X_2)$ :

$$P(X_2 | X_1) \neq P_1(X_2 | \text{do}(X_1)) \neq P_2(X_2 | \text{do}(X_1)) \neq P(X_2 | X_1).$$

This shows that the interventional distribution  $P(X_2 | \text{do}(X_1))$  is not identifiable (and thus not trackable) from  $P(X_1, X_2)$  and  $G$ . However, it is trackable up to oracle choices from  $P(X_1, X_2)$  and  $G$ . It would thus be trackable (and identifiable) for discrete CBNs with strictly positive mass functions, e.g. here if we also put positive mass on  $P(X_1 = c) > 0$ .

**Notation 5.3.6.** Let  $G = (J, V, E, L)$  be a CADMG,  $<$  a topological order of  $G$  and let  $v \in C \subseteq V$ . We then put:

$$\begin{aligned} \text{Anc}^{[C]}(v) &:= \text{Anc}^{G_{\text{do}(C^c)}}(v) \cap C, \\ \text{Pred}_{<}^{[C]}(v) &:= \text{Pred}_{<}^{G_{\text{do}(C^c)}}(v) \cap C &= \text{Pred}_{<}^G(v) \cap C, \\ \text{Dist}^{[C]}(v) &:= \text{Dist}^{G_{\text{do}(C^c)}}(v) \cap C &= \text{Dist}^{G_{\text{do}(C^c)}}(v), \\ \text{Dist}_{<}^{[C]}(v) &:= \text{Dist}^{[C]}(v) \cap \text{Pred}_{<}^{[C]}(v), \\ \mathcal{D}[C] &:= \left\{ \text{Dist}^{[C]}(v) \mid v \in C \right\}. \end{aligned}$$

Similarly, if we use the subscript  $\leq$  we then also include  $v$ . Also note that the dependence on  $G$  in the above constructions is implicit.

**Notation 5.3.7** (The key interventional Markov kernels). Let  $M$  be an  $L$ -CBN with with observable CADMG  $G = (J, V, E, L)$  and  $C \subseteq V$  any subset. We will abbreviate:

$$\mathcal{Q}[C] := P(X_C | \text{do}(X_{J \cup V \setminus C})) = P(X_C | \text{do}(X_{\overline{(J \cup V) \setminus (\text{Pa}^G(C) \cup C)}}), X_{\text{Pa}^G(C) \setminus C})),$$

where the latter identification comes from Lemma 5.3.8 (using the global Markov property), which is “surely” determined, and not just up to some null-set.

Note that we have the corner cases:

$$\mathcal{Q}[V] = P(X_V | \text{do}(X_J)), \quad \mathcal{Q}[\emptyset] = \delta_*.$$

**Lemma 5.3.8.** *Let  $M$  be an  $L$ -CBN with with observable CADMG  $G = (J, V, E, L)$  and  $C \subseteq V$  any subset. Then we have the identification:*

$$P(X_C | \text{do}(X_{J \cup V \setminus C})) = P(X_C | \text{do}(X_{\overline{(J \cup V) \setminus (\text{Pa}^G(C) \cup C)}}), X_{\text{Pa}^G(C) \setminus C})).$$

*Proof.* This follows from the global Markov property with:

$$C \underset{G_{\text{do}(I_{V \setminus (C \cup \text{Pa}^G(C))}, \text{Pa}^G(C) \setminus C)}}{\perp^d} I_{V \setminus (C \cup \text{Pa}^G(C))} | \text{Pa}^G(C) \setminus C.$$

To elaborate the latter, let  $P := \text{Pa}^G(C) \setminus C$  and  $W := V \setminus (C \cup \text{Pa}^G(C))$  and:

$$V' := V \setminus P = C \dot{\cup} W, \quad J' := (J \setminus P) \dot{\cup} P, \quad G' := G_{\text{do}(I_W, P)}.$$

Now consider a walk from a node  $c \in C$  to a node  $j \in J' \dot{\cup} I_W$  in  $G'$ :

$$\pi : \quad c \ast \ast \dots \ast \ast j.$$

If  $j \in P$  then the walk is blocked by  $P$  at the endnode  $j \in P$ . So lets assume the case  $j \notin P$ . Then the walk is of the form:

$$\pi : \quad c \ast \ast \dots \ast \ast w \leftarrow j,$$

with a  $w \in W$ . So we can write it further as:

$$\pi : \quad c = c_0 \ast \ast c_1 \ast \ast \dots c_k \ast \ast v \ast \ast \dots \ast \ast w \leftarrow j,$$

with  $c_0, \dots, c_k \in C$  for some  $k \geq 0$ , and  $v \notin C$ , the first occurring node not in  $C$  (on  $\pi$  from the left). Note that  $v = w$  is possible. If the edge  $c_k \ast \ast v$  is of the form  $c_k \leftarrow v$  then  $v \in P$  and the walk is blocked by  $P$  at the non-collider  $v$ . So we can assume the case where the edge is of the form  $c_k \ast \ast v$ . This means that on the subwalk  $c_k \ast \ast v \ast \ast \dots \ast \ast w \leftarrow j$  we must have at least one collider. This collider is then blocked by  $P$  as no collider can be an ancestor of a node in  $P$  inside  $G'$ , because  $P$  consists only of input nodes of  $G'$ .

This shows the claim:

$$C \underset{G'}{\perp^d} I_W | P.$$

The rest then follows from the global Markov property.  $\square$

### 5.3.2. The Interventional Ordered Local Markov Property

One of the ingredient for the ID-algorithm is the ability to track (up to oracle choices) the interventional Markov kernel  $\mathcal{Q}[D]$  for districts  $D$  of  $G$  from  $\mathcal{Q}[V]$ . The key ingredient to achieve this is the *interventional ordered local Markov property*, which provides us with certain well-behaved Markov kernels that appear in factorizations of both  $\mathcal{Q}[V]$  and  $\mathcal{Q}[D]$ .

**Definition 5.3.9** (The preceding Markov blanket of a node). *Let  $G = (J, V, E, L)$  be a CADMG and  $<$  a topological order. For  $v \in V$  we make the following abbreviations:*

$$\begin{aligned} G_{<}(v) &:= \text{Pred}_{<}^G(v), \\ \text{Di}_{<}^G(v) &:= \text{Dist}^{G_{<}(v)}(v), \\ \text{Di}_{<}^G(v) &:= \text{Di}_{<}^G(v) \setminus \{v\}, \\ \text{PaD}_{<}^G(v) &:= \text{Pa}^G(\text{Di}_{<}^G(v)) \setminus \text{Di}_{<}^G(v), \\ \text{Mb}_{<}^G(v) &:= \text{PaD}_{<}^G(v) \dot{\cup} \text{Di}_{<}^G(v) \\ &= \text{Pa}^G(\text{Dist}^{G_{<}(v)}(v)) \cup \text{Dist}^{G_{<}(v)}(v) \setminus \{v\}. \end{aligned}$$

We call  $\text{Di}_{<}^G(v)$  the preceding district and  $\text{Mb}_{<}^G(v)$  the preceding Markov blanket of  $v$  in  $G$  w.r.t.  $<$ . Note that this definition depends on the topological order  $<$  and that we have the inclusions:

$$\text{Di}_{<}^G(v) \subseteq \text{Mb}_{<}^G(v) \subseteq \text{Pred}_{<}^G(v).$$

**Proposition 5.3.10** (Interventional ordered local Markov property). *Let  $M$  be an L-CBN with with observable CADMG  $G = (J, V, E, L)$  and a fixed topological order  $<$ . Then for every  $v \in V$  we have the conditional independence:*

$$X_v \underset{P(X_V | \text{do}(X_{I_{V \setminus \text{Di}_{<}^G(v)}}, X_J))}{\perp\!\!\!\perp} X_{\text{Pred}_{<}^G(v)} \mid X_{\text{Mb}_{<}^G(v)}.$$

In particular, there exists a Markov kernel, denoted by:

$$Q(X_v | X_{\text{Mb}_{<}^G(v)}) \quad \text{or} \quad Q(X_v | X_{\text{Di}_{<}^G(v)}, \text{do}(X_{\text{PaD}_{<}^G(v)})),$$

that simultaneously is a version of:

$$P(X_v | X_{\text{Pred}_{<}^G(v)}, \text{do}(X_J)) \quad \text{and} \quad P(X_v | X_{\text{Pred}_{<}^G(v)}, \text{do}(X_{J \cup V \setminus D})),$$

for every subset  $D \subseteq V$  with  $\text{Di}_{<}^G(v) \subseteq D$ , e.g.  $D = \text{Dist}^G(v)$ .

*Proof.* This follows from the global Markov property, Theorem 4.2.1, together with the d-separation statement:

$$\{v\} \underset{G_{\text{do}(I_{V \setminus \text{Di}_{<}^G(v)}}}{\perp^d} \text{Pred}_{<}^G(v) \mid \text{Mb}_{<}^G(v).$$

See Lemma 5.3.12 and Lemma 5.3.13. □

## Proofs - The Interventional Ordered Local Markov Property

For the next two Lemmata we introduce some shorter notations:

**Notation 5.3.11.** Let  $G = (J, V, E, L)$  be a CADMG,  $<$  a topological order for  $G$  and  $v \in V$  a fixed node. For a subset  $W \subseteq J \cup V$  we abbreviate:

$$W_{<} := \{w \in W \mid w < v\}, \quad W_{\leq} := \{w \in W \mid w < v \vee w = v\},$$

and  $W_{>}$  and  $W_{\geq}$  accordingly.

The next Lemma is the graphical center piece that makes the ID-algorithm possible. It could be called the graphical version of an “interventional ordered local Markov property”, whose distributional counterpart is stated in the Lemma after.

**Lemma 5.3.12.** Let  $G = (J, V, E, L)$  be a CADMG,  $<$  a topological order for  $G$  and  $v \in V$  a fixed node. Let  $D \subseteq V$  be a subset such that  $v \in D$  and  $\text{Dist}^{G_{\leq}}(D_{\leq}) \subseteq D$ , where  $G_{\leq} := \text{Pred}_{\leq}^G(v)$  is the ancestral subgraph of predecessors of  $v$  in  $G$ , e.g.  $D = \text{Dist}^G(v)$  or  $D = \text{Dist}^{G_{\leq}}(v)$ . Then we have the  $d$ -separation statement:

$$\{v\} \underset{G_{\text{do}(I_{V \setminus D_{\leq}})}}{\perp^d} V_{<} \mid \text{Pa}^G(D_{\leq}) \cup D_{<}.$$

*Proof.* We abbreviate:

$$\tilde{G} := G_{\text{do}(I_{V \setminus D_{\leq}})}, \quad F := \text{Pa}^G(D_{\leq}) \cup D_{<}.$$

Assume the contrary to the claim and let  $\pi$  be a shortest path from  $v$  to a node  $w \in J \cup I_{V \setminus D} \cup V_{<}$  in the graph  $\tilde{G}$ . It is clear that  $w \neq v$ .

If  $w \in F$  then  $\pi$  is blocked by  $F$  at the endnode  $w$ . So we can assume that  $w \notin F$ , in particular,  $w \notin D_{\leq}$ . So there exist  $v_0, \dots, v_k \in D_{\leq}$  for some  $k \geq 0$  and  $\tilde{w} \notin D_{\leq}$  such that  $\pi$  is of the form:

$$\pi : \quad v = v_0 \ast \ast \dots \ast \ast v_k \ast \ast \tilde{w} \ast \ast \dots \ast \ast w,$$

where  $\tilde{w}$  is the first node from the left that is not in  $D_{\leq}$ , which exists since  $w \notin D_{\leq}$ . Since  $v_k \in D_{\leq}$  it is clear that  $\tilde{w} \notin I_{V \setminus D_{\leq}}$ . So  $\tilde{w} \in J \cup V \setminus D_{\leq}$ .

If the edge  $v_k \ast \ast \tilde{w}$  is of the form  $v_k \leftarrow \tilde{w}$  then  $\tilde{w} \in \text{Pa}^G(D_{\leq})$ . So in this case  $\pi$  is blocked at the non-collider  $\tilde{w}$  by  $F$ .

So we can assume the case  $v_k \ast \ast \tilde{w}$ . This implies that  $\tilde{w}$  cannot lie in the set of input nodes of  $\tilde{G}$ , which implies that  $\tilde{w} \notin J \cup I_{V \setminus D_{\leq}}$ . With this we then get that  $\tilde{w} \in V \setminus D_{\leq}$ .

Assume the case that  $\tilde{w} \in V_{>}$  and  $\pi$  is of the form:

$$\pi : \quad v = v_0 \ast \ast \dots \ast \ast v_k \ast \ast \tilde{w} = \tilde{w}_0 \ast \ast \dots \ast \ast \tilde{w}_m = w,$$

where the subwalk  $\tilde{w}_0 \ast \ast \dots \ast \ast \tilde{w}_m$  has no colliders. Because of the edge  $v_k \ast \ast \tilde{w}$  this subwalk necessarily is directed to the right and we get:

$$\pi : \quad v = v_0 \ast \ast \dots \ast \ast v_k \ast \ast \tilde{w} = \tilde{w}_0 \rightarrow \dots \rightarrow \tilde{w}_m = w.$$

Then the endnode  $w$  is not an input node,  $w \notin J \cup I_{V \setminus D_{\leq}}$ , and  $v < \tilde{w} \leq \tilde{w}_m = w$ , thus  $w \in V_{>}$ . But this is a contradiction to  $w \in J \cup I_{V \setminus D} \cup V_{<}$ . So this case cannot occur.

Now assume the case that  $\tilde{w} \in V_{>}$  and  $\pi$  is of the form:

$$\pi : \quad v = v_0 \ast \ast \dots \ast \ast v_k \ast \rightarrow \tilde{w} = \tilde{w}_0 \rightarrow \dots \rightarrow \tilde{w}_m \leftarrow \ast \dots \ast \ast w,$$

for some  $m \geq 0$ , with a directed subwalk  $\tilde{w}_0 \rightarrow \dots \rightarrow \tilde{w}_m$ , where  $\tilde{w}_m$  is the first node after  $\tilde{w}$  where a collider occurs, which could be  $\tilde{w}$  itself. Again we have  $v < \tilde{w} \leq \tilde{w}_m$  and thus  $\tilde{w}_m \in V_{>}$ . This implies that  $\tilde{w}_m \notin \text{Anc}^G(F)$ , since  $\text{Anc}^G(F) \subseteq J \cup V_{<}$ . So  $\pi$  is blocked by  $F$  at the collider  $\tilde{w}_m$ .

So we are left with the cases  $v_k \ast \rightarrow \tilde{w}$  and  $\tilde{w} \in V_{<} \setminus D_{\leq}$ .

Now consider the case of a directed edge  $v_k \rightarrow \tilde{w}$  and  $\tilde{w} \in V_{<} \setminus D_{\leq}$ . If  $v_k \neq v$  then  $\pi$  is blocked at the non-collider  $v_k \in D_{<}$  by  $F$ . So we can assume that  $v_k = v$ . This implies  $v < \tilde{w}$  and thus  $\tilde{w} \in V_{>}$ , which contradicts  $\tilde{w} \in V_{<}$ . So this cannot occur.

Now consider the case of a bidirected edge  $v_k \leftrightarrow \tilde{w}$  and  $\tilde{w} \in V_{<} \setminus D_{\leq}$ . Then both nodes  $v_k, \tilde{w} \in G_{\leq}$  and  $\tilde{w} \in \text{Dist}^{G_{\leq}}(v_k)$ . By the assumption of this Lemma we have:

$$\tilde{w} \in \text{Dist}^{G_{\leq}}(v_k) \subseteq \text{Dist}^{G_{\leq}}(D_{\leq}) \subseteq D \cap G_{\leq} = D_{\leq}.$$

So  $\tilde{w} \in D_{\leq}$ , which contradicts  $\tilde{w} \notin D_{\leq}$ . So this case cannot occur.

So we have shown that in all cases that can occur the path  $\pi$  in  $\tilde{G}$  is blocked by  $F$ . This shows the claim.  $\square$

The last Lemma allows us to use the global Markov property for the existence of special Markov kernels that are of importance for the ID-algorithm:

**Lemma 5.3.13** (Interventional ordered local Markov property). *Let  $M$  be an L-CBN with with observable CADMG  $G = (J, V, E, L)$  and a fixed topological order  $<$  and fixed  $v \in V$ . Let  $G_{\leq} := \text{Pred}_{\leq}^G(v)$  be the ancestral sub-CADMG of predecessors of  $v$  in  $G$  and let:*

$$D_{\leq} := \text{Dist}^{G_{\leq}}(v), \quad F := \text{Pa}^G(D_{\leq}) \cup D_{\leq} \setminus \{v\}.$$

*Then we have the conditional independence:*

$$X_v \perp\!\!\!\perp_{P(X_{V_{\leq}} | \text{do}(X_{I_{V \setminus D_{\leq}}}, X_J))} X_{V_{<}} | X_F.$$

*In particular, there exists a Markov kernel:*

$$Q(X_v | X_F)$$

*that simultaneously is a version of:*

$$P(X_v | X_H, \text{do}(X_{J \cup V \setminus S})),$$

*for every subsets  $H, S \subseteq V$  such that  $D_{\leq} \subseteq S \subseteq V$  and  $F \cap S_{<} \subseteq H \subseteq V_{<}$ . Note that this includes the corner cases:*

1.  $P(X_v | X_{V_{<}}, \text{do}(X_J))$ ,



2.  $P(X_v | X_{F \cap V_{<}}, \text{do}(X_J))$ ,
3.  $P(X_v | X_{S_{<}}, \text{do}(X_{J \cup V \setminus S}))$ ,
4.  $P(X_v | X_{F \cap S_{<}}, \text{do}(X_{J \cup V \setminus S}))$ ,
5.  $P(X_v | X_{D_{\leq}}, \text{do}(X_{J \cup V \setminus D_{\leq}}))$ .

*Proof.* By Lemma 5.3.12 we have the d-separation:

$$\{v\} \underset{G_{\text{do}(I_{V \setminus D_{\leq}})}}{\perp^d} V_{<} | F.$$

By the global Markov property, Theorem 4.2.1, we get:

$$X_v \underset{P(X_{V_{\leq}} | \text{do}(X_{I_{V \setminus D_{\leq}}}, X_J))}{\perp\!\!\!\perp} X_{V_{<}} | X_F.$$

So there exists a Markov kernel  $Q(X_v | X_F)$  such that:

$$P(X_{V_{\leq}} | \text{do}(X_{I_{V \setminus D_{\leq}}}, X_J)) = Q(X_v | X_F) \otimes P(X_{V_{<}} | \text{do}(X_{I_{V \setminus D_{\leq}}}, X_J)). \quad (\#)$$

For a subset  $S \subseteq V$  with  $D_{\leq} \subseteq S$  we have:

$$V \setminus D_{\leq} = (V \setminus S) \dot{\cup} (S \setminus D_{\leq}).$$

By putting  $X_{I_{S \setminus D_{\leq}}} = \star$  and  $X_{I_{V \setminus S}} = x_{V \setminus S}$  we get:

$$P(X_{V_{\leq}} | \text{do}(X_{V \setminus S}, X_J)) = Q(X_v | X_F) \otimes P(X_{V_{<}} | \text{do}(X_{V \setminus S}, X_J)).$$

Now consider another subset  $H \subseteq V_{<}$  with  $F \cap S_{<} \subseteq H$ . Then marginalizing out  $X_{V_{\leq} \setminus (\{v\} \cup H)}$  gives us:

$$P(X_v, X_H | \text{do}(X_{J \cup V \setminus S})) = Q(X_v | X_F) \otimes P(X_H | \text{do}(X_{J \cup V \setminus S})).$$

This shows that  $Q(X_v | X_F)$ , simultaneously, is a version of:

$$P(X_v | X_H, \text{do}(X_{J \cup V \setminus S})),$$

for every subsets  $H, S \subseteq V$  such that  $D_{\leq} \subseteq S \subseteq V$  and  $F \cap S_{<} \subseteq H \subseteq V_{<}$ . This shows the claim.  $\square$

### 5.3.3. Ancestral Sets and Districts

**Lemma 5.3.14** (Ancestral subsets are trackable). *Let  $M$  be an  $L$ -CBN with with observable CADMG  $G = (J, V, E, L)$  and  $A \subseteq V$  be a subset such that  $A = \text{Anc}^{[V]}(A)$ . Then we have the equality between the interventional distribution  $\mathcal{Q}[A]$  and the  $A$ -marginal of  $\mathcal{Q}[V]$ :*

$$\mathcal{Q}[A] = P(X_A | \text{do}(X_{J \cup V \setminus A})) = P(X_A | \text{do}(X_J)).$$

*Proof.* By Lemma 5.3.8 we only have to show the right identity:

$$\mathcal{Q}[A] = P(X_A | \text{do}(X_{J \cup V \setminus A})) \stackrel{!}{=} P(X_A | \text{do}(X_J)).$$

The latter follows again from the global Markov property and the d-separation:

$$A \underset{G_{\text{do}(I_{V \setminus A})}}{\perp^d} I_{V \setminus A} | J.$$

This d-separation holds true since every walk from a node  $v \in A$  to a node  $j \in J \cup I_{V \setminus A}$  is either blocked by  $J$  as the endnode  $j \in J$  or is of the form:

$$\pi : \quad v \ast\ast a_1 \ast\ast \dots \ast\ast a_k \ast\ast w' \ast\ast \dots \ast\ast w \leftarrow j,$$

with a  $w \in V \setminus A$ ,  $j \in I_{V \setminus A}$ ,  $a_1, \dots, a_k \in A$  for some  $k \geq 0$  and  $w' \notin A$ , the first node not in  $A$  on  $\pi$  from the left ( $w' = w$  possible). In case the edge  $a_k \ast\ast w'$  is of the form  $a_k \leftarrow w'$  we have:

$$w' \in \text{Anc}^G(A) \setminus A = \text{Anc}^G(A) \setminus \text{Anc}^{[V]}(A) \subseteq J.$$

So in this case the walk is blocked at the non-collider  $w'$  by  $J$ . So we can consider the case where the edge is of the form  $a_k \ast\ast w'$ . Then the subwalk  $a_k \ast\ast w' \ast\ast \dots \ast\ast w \leftarrow j$  must contain a collider. This collider can not be an ancestor of  $J$ , as  $J$  are the input nodes. So the walk is blocked by  $J$  in all cases.  $\square$

**Remark 5.3.15** (Districts are trackable up to oracle choices). *Let  $M$  be an L-CBN with observable CADMG  $G = (J, V, E, L)$ .*

1. *Since the Markov kernel  $Q(X_v | X_{\text{Mb}_{\leq}^G(v)})$  coming from the interventional ordered local Markov property, see Proposition 5.3.10, is a version of both:*

$$P(X_v | X_{\text{Pred}_{\leq}^{[V]}(v)}, \text{do}(X_J)) \quad \text{and} \quad P(X_v | X_{\text{Pred}_{\leq}^{[D]}(v)}, \text{do}(X_{J \cup V \setminus D})),$$

*for  $D = \text{Dist}^G(v)$ , which are marginal conditionals of the interventional distributions  $\mathcal{Q}[V]$  and  $\mathcal{Q}[D]$ , resp., the  $Q(X_v | X_{\text{Mb}_{\leq}^G(v)})$ 's are trackable from either quantity up to oracle choices.*

2. *We get the following factorization by the chain rule for every  $D \in \mathcal{D}[V]$ :*

$$\begin{aligned} \mathcal{Q}[V] &= \bigotimes_{v \in V}^{\succ} P(X_v | X_{\text{Pred}_{\leq}^{[V]}(v)}, \text{do}(X_J)) &= \bigotimes_{v \in V}^{\succ} Q(X_v | X_{\text{Mb}_{\leq}^G(v)}), \\ \mathcal{Q}[D] &= \bigotimes_{v \in D}^{\succ} P(X_v | X_{\text{Pred}_{\leq}^{[D]}(v)}, \text{do}(X_{J \cup V \setminus D})) &= \bigotimes_{v \in D}^{\succ} Q(X_v | X_{\text{Mb}_{\leq}^G(v)}). \end{aligned}$$

3. *In particular,  $\mathcal{Q}[D]$  is trackable from  $\mathcal{Q}[V]$  up to oracle choices by first determining  $Q(X_v | X_{\text{Mb}_{\leq}^G(v)})$  for  $v \in D$  via marginalization and conditioning and then taking the product of Markov kernels in reverse order of  $\prec$ .*

4. The above seems to give us something like a factorization:

$$\begin{aligned} \mathcal{Q}[V] &= \bigotimes_{v \in V}^{\succ} Q(X_v | X_{\text{Mb}_{\prec}^G(v)}) &= \left[ \bigotimes_{D \in \mathcal{D}[V]}^{\succ} \bigotimes_{v \in D}^{\succ} \right] Q(X_v | X_{\text{Mb}_{\prec}^G(v)}) \\ &= \left[ \bigotimes_{D \in \mathcal{D}[V]}^{\succ} \right] \left( \bigotimes_{v \in D}^{\succ} Q(X_v | X_{\text{Mb}_{\prec}^G(v)}) \right) &= \left[ \bigotimes_{D \in \mathcal{D}[V]}^{\succ} \right] \mathcal{Q}[D], \end{aligned}$$

where the products in brackets are not well-defined in the naive way, as the districts of a CADMG don't need to be topologically ordered. Nonetheless, if we are given a fixed topological order  $\prec$  and the interventional Markov kernels  $\mathcal{Q}[D]$  for every  $D \in \mathcal{D}[V]$  then we can first track  $Q(X_v | X_{\text{Mb}_{\prec}^G(v)})$  from  $\mathcal{Q}[D]$  up to oracle choices for every  $v \in D$  and every  $D \in \mathcal{D}[V]$  and then take the product on the top left.

**Definition 5.3.16.** Let  $M$  be an  $L$ -CBN with observable CADMG  $G = (J, V, E, L)$  and a fixed topological order  $\prec$ . Let  $\mathcal{D} \subseteq \mathcal{D}[V]$  be a set of districts of  $G$ . Then we put:

$$\left[ \bigotimes_{D \in \mathcal{D}}^{\succ} \right] \mathcal{Q}[D] := \bigotimes_{v \in \bigcup_{D \in \mathcal{D}} D}^{\succ} Q(X_v | X_{\text{Mb}_{\prec}^G(v)}),$$

where the product on the right is taken in reverse topological order. Note that for  $G' := G_{\text{do}(D^c)}$  and  $v \in D$  we have  $\text{Mb}_{\prec}^{G'}(v) = \text{Mb}_{\prec}^G(v)$ . So the set  $\text{Mb}_{\prec}^G(v)$  can be determined by the subgraph  $G'$  and  $\prec$  alone.

### 5.3.4. The ID-Algorithm

Now we come to the main part of this section, the *ID-algorithm* for the identification of causal effects, or, more precisely, the trackability up to oracle choices of interventional Markov kernels, from the observable Markov kernel. The main references are [Pea09, GP95, Tia02, TP02, Tia04, SP06b, HV06, HV08, RERS23, FM20].

**Algorithm 5.3.17** (ID-algorithm). Let  $M$  be an  $L$ -CBN with with observable CADMG  $G = (J, V, E, L)$  and a fixed topological order  $\prec$ . Let  $\emptyset \neq B \subseteq V$  and  $C \subseteq J \cup V$  be two disjoint subsets of nodes. We want to query if the interventional Markov kernel  $P(X_B | \text{do}(X_{J \cup C}))$  is trackable up to oracle choices from the observable Markov kernel  $P(X_V | \text{do}(X_J)) = \mathcal{Q}[V]$  and  $G$ .

1. Put  $B^C := \text{Anc}^{G_{\text{do}(C)}}(B) \setminus (J \cup C) \subseteq V$ .

Then  $P(X_B | \text{do}(X_{J \cup C}))$  is the  $B$ -marginal of  $P(X_{B^C} | \text{do}(X_{J \cup C})) = \mathcal{Q}[B^C]$ .

So we are left to determine if we can track  $\mathcal{Q}[B^C]$  from  $\mathcal{Q}[V]$  up to oracle choices.

2. Find the districts  $\mathcal{D}[B^C] = \{S_1, \dots, S_K\}$  and put  $A_{k,0} := V$  for  $k = 1, \dots, K$ .

Note that  $\mathcal{Q}[A_{k,0}] = \mathcal{Q}[V]$  is trivially tracked from  $\mathcal{Q}[V]$ .

3. For each  $k = 1, \dots, K$  repeat the following steps recursively for  $\ell \in \mathbb{N}$ :

a) Take the district in  $A_{k,\ell}$ :

$$D_{k,\ell} := \text{Dist}^{[A_{k,\ell}]}(S_k)$$

We can track  $\mathcal{Q}[D_{k,\ell}]$  from  $\mathcal{Q}[A_{k,\ell}]$  up to oracle choices by Remark 5.3.15.

b) Take the ancestral closure in  $D_{k,\ell}$ :

$$A_{k,\ell+1} := \text{Anc}^{[D_{k,\ell}]}(S_k).$$

We can track  $\mathcal{Q}[A_{k,\ell+1}]$  from  $\mathcal{Q}[D_{k,\ell}]$  via marginalization by Lemma 5.3.14.

c) If  $D_{k,\ell} = A_{k,\ell}$  or  $A_{k,\ell+1} = D_{k,\ell}$  then stop for this  $k$  and put:

$$\check{S}_k := D_{k,\ell}.$$

Otherwise, repeat with:  $\ell \leftarrow \ell + 1$ .

4. When the algorithm has stopped then for every  $k = 1, \dots, K$  we have:

$$\check{S}_k = \text{Anc}^{[\check{S}_k]}(S_k) = \text{Dist}^{[\check{S}_k]}(S_k) \supseteq S_k.$$

Furthermore, we have tracked all  $\mathcal{Q}[\check{S}_k]$ 's recursively from  $\mathcal{Q}[V]$  up to oracle choices.

5. If there is any  $k = 1, \dots, K$  with  $\check{S}_k \neq S_k$  then the ID-algorithm outputs: FAIL.

6. Otherwise, we have for all  $k = 1, \dots, K$  that  $\mathcal{Q}[S_k] = \mathcal{Q}[\check{S}_k]$  and we can track  $\mathcal{Q}[B^C]$  from  $\mathcal{Q}[V]$  up to oracle choices via:

$$\mathcal{Q}[B^C] = \left[ \begin{array}{c} \bigotimes_{S_k \in \mathcal{D}[B^C]} \\ \text{>} \end{array} \right] \mathcal{Q}[S_k],$$

and  $P(X_B | \text{do}(X_{J \cup C}))$  as the  $B$ -marginal thereof.

**Corollary 5.3.18** (Soundness up to oracle choices). *The ID-algorithm 5.3.17 is sound up to oracle choices. This means that if it does not produce FAIL for input  $B, C \subseteq G$  then  $P(X_B | \text{do}(X_{J \cup C}))$  is trackable up to oracle choices from  $P(X_V | \text{do}(X_J))$  and  $G$ .*

*Proof.* This is clear as each step in the ID-algorithm is trackable up to oracle choices. Note that the operations at each step can be formulated by knowing  $G$ ,  $B$  and  $C$  alone without knowing  $M$  in advance.  $\square$

**Theorem 5.3.19** (Soundness up to null-sets). *Let  $G = (J, V, E, L)$  be a CADMG with a fixed topological order  $<$ . Consider the class of  $L$ -CBNs  $M$  with observable CADMG  $G$  such that the following holds:*

1. Each measurable space  $\mathcal{X}_v$  comes equipped with a fixed measure  $\mu_v$  for  $v \in V$ ,

2. for every subset  $D \subseteq V$  the interventional Markov kernel  $\mathcal{Q}[D] = P(X_D | \text{do}(X_{J \cup V \setminus D}))$  is absolute continuous w.r.t. the product measure  $\mu_D := \bigotimes_{v \in D} \mu_v$  and vice versa:

$$\mu_D \ll \mathcal{Q}[D] \ll \mu_D.$$

If the ID-algorithm does not produce FAIL for input  $B, C \subseteq G$ , then  $P(X_B | \text{do}(X_{J \cup C}))$  is “almost-surely” trackable from  $P(X_V | \text{do}(X_J))$  and  $G$  for such CBNs  $M$ , i.e. the Markov kernel that was output by the ID-algorithm equals  $P(X_B | \text{do}(X_{J \cup C}))$  up to a  $\mu_{V \setminus B^C}$ -null set in  $\mathcal{X}_{J \cup V \setminus B^C}$  (see Remark 5.3.20 below).

*Proof.* See Theorem 5.3.31. □

**Remark 5.3.20.** 1. For the almost-sure soundness in Theorem 5.3.19 to hold one implicitly needs/is allowed to make slight relaxations to the ID-algorithm 5.3.17:

Instead of insisting on taking the conditionals  $Q(X_v | X_{\text{Mb}_{\xi'}(v)})$  of  $\mathcal{Q}[D]$  for  $v \in D \subseteq V$  one takes any version of that conditional of  $\mathcal{Q}[D] \otimes \mu_{V \setminus D}$  that is only dependent on predecessors of  $v$ , which will always be possible as  $Q(X_v | X_{\text{Mb}_{\xi'}(v)})$  is an existing such version.

2. The conditions of Theorem 5.3.19 are satisfied for a L-CBNs  $M$  with CDAG  $G^+ = (J, U \dot{\cup} V, E)$  if every Markov kernel  $P_v(X_v | \text{do}(X_{\text{Pa}^{G^+}(v)}))$  has a strictly positive Doob-Radon-Nikodym derivative/density w.r.t. the measure  $\mu_v$  for all  $v \in V$ , see Lemma 5.3.33.

**Theorem 5.3.21** (From almost-sure to sure soundness). *In addition to the conditions in Theorem 5.3.19 assume:*

1.  $\mathcal{X}_v$  is a Polish space for every  $v \in J \cup V$ ,
2.  $\mu_V$  is strictly positive (on non-empty open subsets of  $\mathcal{X}_V$ ),
3. the queried interventional Markov kernel is continuous as a map:

$$P(X_B | \text{do}(X_{J \cup C})) : \mathcal{X}_{J \cup C} \rightarrow \mathcal{P}(\mathcal{X}_B).$$

Then the output  $\hat{P}(X_B | X_{J \cup V \setminus B^C})$  of the ID-algorithm (in the not-FAIL case) can be changed on a  $\mu_{V \setminus B^C}$ -null set in  $\mathcal{X}_{J \cup V \setminus B^C}$  such that it becomes continuous as a map:

$$\hat{P}(X_B | X_{J \cup V \setminus B^C}) : \mathcal{X}_{J \cup V \setminus B^C} \rightarrow \mathcal{P}(\mathcal{X}_B).$$

Every such continuous version of  $\hat{P}(X_B | X_{J \cup V \setminus B^C})$  is then necessarily identical to the interventional Markov kernel  $P(X_B | \text{do}(X_{J \cup C}))$ . These conditions thus allow us to recover from the ambiguity resulting from the null-sets.

*Proof.* The existence of such a null-set is clear, because  $\hat{P}(X_B | X_{J \cup V \setminus B^C})$  is  $\mu_{V \setminus B^C}$ -almost-surely equal to  $P(X_B | \text{do}(X_{J \cup C}))$ , and the latter was assumed to be continuous. The uniqueness follows from Lemma 2.4.23 □

**Remark 5.3.22.** *The statement of Theorem 5.3.21 can be further relaxed by asking for Polish spaces  $\mathcal{X}_v$  only for  $v \in V$ , strict positivity only for  $\mu_{V \setminus B^c}$  and only for the continuity of the maps:*

$$\mathcal{X}_{C \setminus J} \rightarrow \mathcal{P}(\mathcal{X}_B), \quad x_{C \setminus J} \mapsto P(X_B | \text{do}(X_{J \cup C} = (x_J, x_{C \setminus J}))),$$

for every  $x_J \in \mathcal{X}_J$  separately, by then applying the criterion from Theorem 5.3.21 for each partial input  $x_J \in \mathcal{X}_J$  separately.

**Example 5.3.23.** *All stated assumptions of Theorems 5.3.19 and 5.3.21 are satisfied if every Markov kernel of the L-CBN  $M$  is linear Gaussian:*

$$P_v(X_v \in dx_v | \text{do}(X_{\text{Pa}^G+(v)} = x_{\text{Pa}^G+(v)})) = \mathcal{N}(dx_v | \Gamma_v \cdot x_{\text{Pa}^G+(v)} + \gamma_v, \Sigma_v),$$

with transition matrix  $\Gamma_v$ , translation vector  $\gamma_v$  and positive definite covariance matrix  $\Sigma_v \succ 0$  and Lebesgue measures  $\mu_v$ ,  $v \in U \cup V$ .

**Theorem 5.3.24** (Completeness, see [HV08]). *The ID-algorithm is complete.*

More precisely, if the ID-algorithm outputs FAIL for subsets  $B, C \subseteq G = (J, V, E, L)$  then there exist two L-CBNs  $M_1$  and  $M_2$  with the same observable CADMG  $G$ , the same and discrete underlying spaces  $\mathcal{X}_v$  for  $v \in J \cup V$ , and the same observable Markov kernels  $P_1(X_V | \text{do}(X_J)) = P_2(X_V | \text{do}(X_J))$  that have strictly positive mass functions such that:

$$P_1(X_B | \text{do}(X_{J \cup C})) \neq P_2(X_B | \text{do}(X_{J \cup C})).$$

In particular, in case of FAIL,  $P(X_B | \text{do}(X_{J \cup C}))$  is not identifiable from  $P(X_V | \text{do}(X_J))$  and  $G$ .

**Remark 5.3.25** (Identification of conditional causal effects, see [Tia04]). *If we want to know if the conditional interventional Markov kernel  $P(X_A | X_B, \text{do}(X_{J \cup C}))$  is trackable up to oracle choices from  $P(X_V | \text{do}(X_J))$  and  $G$  then we can run the ID-algorithm for  $A \cup B$  and  $C$ . If it does not output FAIL then  $P(X_{A \cup B} | \text{do}(X_{J \cup C}))$  is trackable up to oracle choices from  $P(X_V | \text{do}(X_J))$  and  $G$  and by conditioning on  $X_B$  afterwards so will  $P(X_A | X_B, \text{do}(X_{J \cup C}))$  be.*

However, note that there is a “conditional” version of the ID-algorithm, see [Tia04], that can check if (and conclude that)  $P(X_A | X_B, \text{do}(X_{J \cup C}))$  is trackable up to oracle choices from  $P(X_V | \text{do}(X_J))$  and  $G$  even if the (unconditional) ID-algorithm outputs FAIL for  $P(X_{A \cup B} | \text{do}(X_{J \cup C}))$ .

## Examples

**Example 5.3.26.** *Consider the DAG  $G = (V, E)$  from Figure 10 (a) with  $V = \{v_1, v_2\}$ ,  $E = \{v_1 \rightarrow v_2\}$ . We want to determine if we can identify  $P(X_2 | \text{do}(X_1))$  from  $P(X_1, X_2)$  in case we have a discrete CBN with strictly positive mass function:  $p(x_1, x_2, x_3) > 0$ . For this let  $B := \{v_2\}$  and  $C := \{v_1\}$ . Note that we have the topological order  $v_1 < v_2$ . We then follow the steps of the ID-algorithm:*

1.  $\mathcal{Q}[V](x_1, x_2) := p(x_1, x_2)$ .
2.  $B^C = \text{Anc}^{G_{\text{do}(C)}}(B) \setminus C = \{v_2\}$ .
3.  $\mathcal{D}[B^C] = \{S = \{v_2\}\}$ . So we compute:
  - a)  $D_0 = \text{Dist}^{[V]}(S) = \{v_2\}$ . Compute:

$$\mathcal{Q}[D_0](x_2|x_1) = q(x_2|x_1) = \frac{\mathcal{Q}[V](x_1, x_2)}{\mathcal{Q}[V](x_1)} = p(x_2|x_1).$$

- b)  $A_1 = \text{Anc}^{[D_0]}(S) = \{v_2\} = D_0$ , thus  $\check{S} = D_0 = \{v_2\} = S$ . Compute:

$$\mathcal{Q}[\check{S}](x_2|x_1) = \mathcal{Q}[D_0](x_2|x_1) = p(x_2|x_1).$$

4. Since  $\check{S} = S = \{v_2\}$  we can compute:

$$\mathcal{Q}[B^C](x_2|x_1) = \mathcal{Q}[\check{S}](x_2|x_1) = p(x_2|x_1).$$

So we can identify  $P(X_2|\text{do}(X_1))$  from  $P(X_1, X_2)$  as the conditional  $P(X_2|X_1)$  via the mass function from above.

**Example 5.3.27.** Consider the ADMG  $G = (V, E, L)$  from Figure 10 (b) with  $V = \{v_1, v_2\}$ ,  $E = \{v_1 \rightarrow v_2\}$  and  $L = \{v_1 \leftrightarrow v_2\}$ . We want to determine if we can identify  $P(X_2|\text{do}(X_1))$  from  $P(X_1, X_2)$  in case we have a discrete CBN with strictly positive mass function:  $p(x_1, x_2, x_3) > 0$ . For this let  $B := \{v_2\}$  and  $C := \{v_1\}$ . Note that we have the topological order  $v_1 < v_2$ . We then follow the steps of the ID-algorithm:

1.  $\mathcal{Q}[V](x_1, x_2) := p(x_1, x_2)$ .
2.  $B^C = \text{Anc}^{G_{\text{do}(C)}}(B) \setminus C = \{v_2\}$ .
3.  $\mathcal{D}[B^C] = \{S = \{v_2\}\}$ . So we compute:
  - a)  $D_0 = \text{Dist}^{[V]}(S) = \{v_1, v_2\}$ . Compute:

$$\begin{aligned} q(x_1) &= \mathcal{Q}[V](x_1) &&= p(x_1), \\ q(x_2|x_1) &= \frac{\mathcal{Q}[V](x_1, x_2)}{\mathcal{Q}[V](x_1)} &&= p(x_2|x_1), \\ \mathcal{Q}[D_0](x_1, x_2) &= q(x_2|x_1) \cdot q(x_1) &&= p(x_1, x_2). \end{aligned}$$

- b)  $A_1 = \text{Anc}^{[D_0]}(S) = \{v_1, v_2\} = D_0$ , thus  $\check{S} = D_0 = \{v_1, v_2\}$ . Compute:

$$\mathcal{Q}[\check{S}](x_1, x_2) = \mathcal{Q}[D_0](x_1, x_2) = p(x_1, x_2).$$

4. Since  $\check{S} = \{v_1, v_2\} \neq \{v_2\} = S$  the ID-algorithms outputs: *FAIL*.

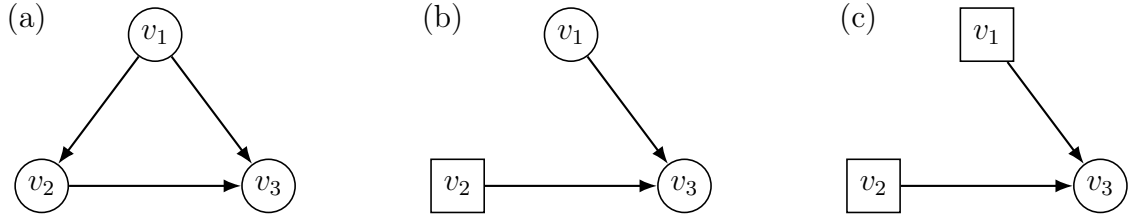


Figure 11: A DAG with three nodes and its intervened graphs.

So we can not identify  $P(X_2 | \text{do}(X_1))$  from  $P(X_1, X_2)$  and  $G$ .

**Example 5.3.28.** Consider the DAG  $G = (V, E, L)$  from Figure 11 with  $V = \{v_1, v_2, v_3\}$ ,  $E = \{v_1 \rightarrow v_2, v_1 \rightarrow v_3, v_2 \rightarrow v_3\}$ . We want to determine if we can identify  $P(X_3 | \text{do}(X_2))$  from  $P(X_1, X_2, X_3)$  in case we have a discrete CBN with strictly positive mass function:  $p(x_1, x_2, x_3) > 0$ . For this let  $B := \{v_3\}$  and  $C := \{v_2\}$ . Note that we have the topological order  $v_1 < v_2 < v_3$ . We then follow the steps of the ID-algorithm:

1.  $\mathcal{Q}[V](x_1, x_2, x_3) := p(x_1, x_2, x_3)$ .
2.  $B^C = \text{Anc}^{G_{\text{do}(C)}}(B) \setminus C = \{v_1, v_3\}$ .
3.  $\mathcal{D}[B^C] = \{S_1 = \{v_1\}, S_2 = \{v_3\}\}$ , see Figure 11 (b).
4. For  $S_1 = \{v_1\}$ :

- a)  $D_{1,0} = \text{Dist}^{[V]}(S_1) = \{v_1\}$ . Compute:

$$\mathcal{Q}[D_{1,0}](x_1) = q(x_1) = \mathcal{Q}[V](x_1) = p(x_1).$$

- b)  $A_{1,1} = \text{Anc}^{[D_{1,0}]}(S_1) = \{v_1\} = D_{1,0}$ , thus  $\check{S}_1 = D_{1,0}$ . Compute:

$$\mathcal{Q}[\check{S}_1](x_1) = \mathcal{Q}[D_{1,0}](x_1) = p(x_1).$$

- c)  $S_1 = \{v_1\} = \check{S}_1$ . Compute:

$$\mathcal{Q}[S_1](x_1) = \mathcal{Q}[\check{S}_1](x_1) = p(x_1).$$

5. For  $S_2 = \{v_3\}$ :

- a)  $D_{2,0} = \text{Dist}^{[V]}(S_2) = \{v_3\}$ . Compute:

$$\mathcal{Q}[D_{2,0}](x_3 | x_1, x_2) = q(x_3 | x_1, x_2) = \frac{\mathcal{Q}[V](x_1, x_2, x_3)}{\mathcal{Q}[V](x_1, x_2)} = p(x_3 | x_1, x_2).$$

- b)  $A_{2,1} = \text{Anc}^{[D_{2,0}]}(S_2) = \{v_3\} = D_{2,0}$ . thus  $\check{S}_2 = D_{2,0}$ . Compute:

$$\mathcal{Q}[\check{S}_2](x_3 | x_1, x_2) = \mathcal{Q}[D_{2,0}](x_3 | x_1, x_2) = p(x_3 | x_1, x_2).$$



c)  $S_2 = \{v_3\} = \check{S}_2$ . Compute:

$$\mathcal{Q}[S_2](x_3|x_1, x_2) = \mathcal{Q}[\check{S}_2](x_3|x_1, x_2) = p(x_3|x_1, x_2).$$

6. Since both  $\check{S}_1 = S_1 = \{v_1\}$  and  $\check{S}_2 = S_2 = \{v_3\}$  we can compute:

$$\begin{aligned} \mathcal{Q}[B^C](x_1, x_3|x_2) &= \mathcal{Q}[S_2](x_3|x_1, x_2) \cdot \mathcal{Q}[S_1](x_1) \\ &= p(x_3|x_1, x_2) \cdot p(x_1), \\ p(x_3|\text{do}(x_2)) &= \sum_{x_1} \mathcal{Q}[B^C](x_1, x_3|x_2) \\ &= \sum_{x_1} p(x_3|x_1, x_2) \cdot p(x_1). \end{aligned}$$

So we can identify  $P(X_3|\text{do}(X_2))$  from  $P(X_1, X_2, X_3)$  via the mass function from above.

**Example 5.3.29** (Counter example when mass functions are not strictly positive). Consider the DAG  $G = (V, E, L)$  from Figure 11 with  $V = \{v_1, v_2, v_3\}$ ,  $E = \{v_1 \rightarrow v_2, v_1 \rightarrow v_3, v_2 \rightarrow v_3\}$ . We want to determine if we can identify  $P(X_3|\text{do}(X_2))$  from  $P(X_1, X_2, X_3)$  in case we do NOT have a strictly positive mass function:  $p(x_1, x_2, x_3) > 0$ . We assume  $\mathcal{X}_1 = \mathcal{X}_2 = \mathcal{X}_3 = \{0, 1\}$ .

$$\begin{aligned} p(x_1 = 1) &:= \frac{1}{2}, & p(x_2|\text{do}(x_1)) &:= \delta_{x_1}(x_2), \\ p(x_3 = 1|\text{do}(x_1 = 0, x_2 = 0)) &:= \frac{1}{4}, & p(x_3 = 1|\text{do}(x_1 = 1, x_2 = 0)) &:= \frac{3}{4}, \\ p(x_3 = 1|\text{do}(x_1 = 0, x_2 = 1)) &:= \frac{1}{8}, & p(x_3 = 1|\text{do}(x_1 = 1, x_2 = 1)) &:= \frac{3}{8}. \end{aligned}$$

Then we get for the  $(X_1, X_2)$ -marginal of the observational distribution  $P(X_1, X_2) = P(X_2|\text{do}(X_1)) \otimes P(X_1)$  the mass functions:

$$\begin{aligned} p(x_1 = 0, x_2 = 0) &= p(x_1 = 1, x_2 = 1) = \frac{1}{2}, \\ p(x_1 = 1, x_2 = 0) &= p(x_1 = 0, x_2 = 1) = 0. \end{aligned}$$

Note that this shows that  $P(X_1, X_2)$  and thus  $P(X_1, X_2, X_3)$  do not have a strictly positive mass functions. With this a valid conditional for the observational joint distribution  $P(X_1, X_2, X_3)$  conditioned on  $(X_1, X_2)$  is:

$$\begin{aligned} p(x_3 = 1|x_1 = 0, x_2 = 0) &:= \frac{1}{4}, & p(x_3 = 1|x_1 = 1, x_2 = 0) &:= \frac{3}{8}, \\ p(x_3 = 1|x_1 = 0, x_2 = 1) &:= \frac{5}{8}, & p(x_3 = 1|x_1 = 1, x_2 = 1) &:= \frac{3}{8}. \end{aligned}$$

The interventional distribution is given by:

$$\begin{aligned}
p(x_3 | \text{do}(x_2)) &= \sum_{x_1} p(x_3 | \text{do}(x_1, x_2)) \cdot p(x_1), \\
p(x_3 = 1 | \text{do}(x_2 = 0)) &= \frac{1}{2} (p(x_3 = 1 | \text{do}(x_1 = 0, x_2 = 0)) + p(x_3 = 1 | \text{do}(x_1 = 1, x_2 = 0))) \\
&= \frac{1}{2} \left( \frac{1}{4} + \frac{3}{4} \right) = \frac{1}{2}, \\
p(x_3 = 1 | \text{do}(x_2 = 1)) &= \frac{1}{2} (p(x_3 = 1 | \text{do}(x_1 = 0, x_2 = 1)) + p(x_3 = 1 | \text{do}(x_1 = 1, x_2 = 1))) \\
&= \frac{1}{2} \left( \frac{1}{8} + \frac{3}{8} \right) = \frac{1}{4}.
\end{aligned}$$

On the other hand, using the other conditional mass functions instead, gives us:

$$\begin{aligned}
\hat{p}(x_3 | x_2) &:= \sum_{x_1} p(x_3 | x_1, x_2) \cdot p(x_1), \\
\hat{p}(x_3 = 1 | x_2 = 0) &= \frac{1}{2} (p(x_3 = 1 | x_1 = 0, x_2 = 0) + p(x_3 = 1 | x_1 = 1, x_2 = 0)) \\
&= \frac{1}{2} \left( \frac{1}{4} + \frac{1}{2} \right) = \frac{3}{8} \\
&\neq \frac{1}{2} = p(x_3 = 1 | \text{do}(x_2 = 0)), \\
\hat{p}(x_3 = 1 | x_2 = 1) &= \frac{1}{2} (p(x_3 = 1 | x_1 = 0, x_2 = 1) + p(x_3 = 1 | x_1 = 1, x_2 = 1)) \\
&= \frac{1}{2} \left( \frac{5}{8} + \frac{1}{8} \right) = \frac{3}{8} \\
&\neq \frac{1}{4} = p(x_3 = 1 | \text{do}(x_2 = 1)).
\end{aligned}$$

This shows that the “surrogate” Markov kernel  $\hat{P}(X_3 | X_2) := P(X_3 | X_1, X_2) \circ P(X_1)$ , which would be proposed by both the ID-algorithm and the backdoor criterion, is NOT equal to the interventional Markov kernel  $P(X_3 | \text{do}(X_2)) = P(X_3 | \text{do}(X_1, X_2)) \circ P(X_1)$ , not even  $P(X_2)$ -almost-surely.

**Example 5.3.30.** Consider the ADMG  $G = (V, E, L)$  from Figure 12 with  $V = \{v_1, v_2, v_3\}$ ,  $E = \{v_1 \rightarrow v_2, v_2 \rightarrow v_3\}$ ,  $J = \emptyset$  and  $L = \{v_1 \leftrightarrow v_3\}$ . We want to determine if we can identify  $P(X_3 | \text{do}(X_1))$  from  $P(X_1, X_2, X_3)$  in case we have a discrete CBN with strictly positive mass function:  $p(x_1, x_2, x_3) > 0$ . For this let  $B := \{v_3\}$  and  $C := \{v_1\}$ . Note that we have the topological order  $v_1 < v_2 < v_3$ . We then follow the steps of the ID-algorithm:

1.  $\mathcal{Q}[V](x_1, x_2, x_3) := p(x_1, x_2, x_3)$ .
2.  $B^C = \text{Anc}^{G_{\text{do}(C)}}(B) \setminus C = \{v_2, v_3\}$ .

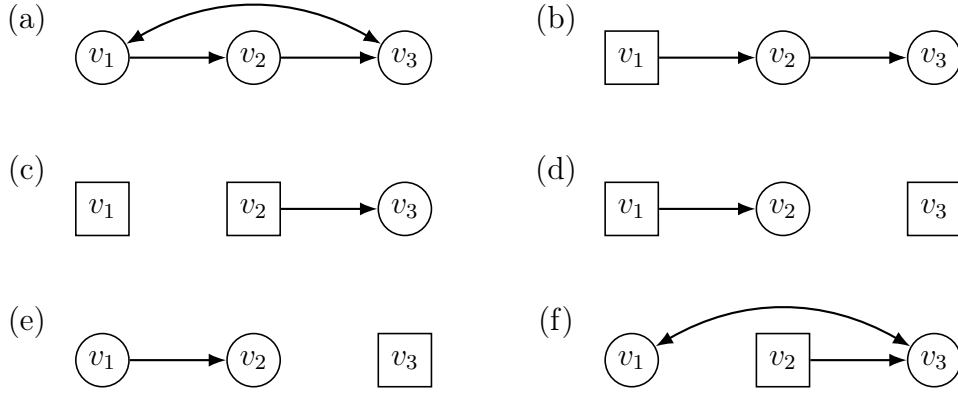


Figure 12: An ADMG and its mutilations, corresponding to the interventional Markov kernels: (a)  $\mathcal{Q}[\{v_1, v_2, v_3\}]$ , (b)  $\mathcal{Q}[\{v_2, v_3\}]$ , (c)  $\mathcal{Q}[\{v_3\}]$ , (d)  $\mathcal{Q}[\{v_2\}]$ , (e)  $\mathcal{Q}[\{v_1, v_2\}]$ , (f)  $\mathcal{Q}[\{v_2, v_3\}]$ .

3.  $\mathcal{D}[B^C] = \{S_1 = \{v_3\}, S_2 = \{v_2\}\}$ , see Figure 12 (b).

4. For  $S_1 = \{v_3\}$ :

a)  $D_{1,0} = \text{Dist}^{[V]}(S_1) = \{v_1, v_3\}$ , see Figure 12 (a), (f). Compute:

$$\begin{aligned} q(x_1) &= \mathcal{Q}[V](x_1) &&= p(x_1), \\ q(x_3|x_1, x_2) &= \frac{\mathcal{Q}[V](x_1, x_2, x_3)}{\mathcal{Q}[V](x_1, x_2)} &&= p(x_3|x_1, x_2), \\ \mathcal{Q}[D_{1,0}](x_1, x_3|x_2) &= q(x_3|x_1, x_2) \cdot q(x_1) &&= p(x_3|x_1, x_2) \cdot p(x_1). \end{aligned}$$

b)  $A_{1,1} = \text{Anc}^{[D_{1,0}]}(S_1) = \{v_3\}$ . see Figure 12 (f), (c). Compute:

$$\mathcal{Q}[A_{1,1}](x_3|x_2) = \sum_{x_1} \mathcal{Q}[D_{1,0}](x_1, x_3|x_2) = \sum_{x_1} p(x_3|x_1, x_2) \cdot p(x_1).$$

c)  $D_{1,1} = \text{Dist}^{[A_{1,1}]}(S_1) = \{v_3\} = A_{1,1}$ , thus  $\check{S}_1 = A_{1,1}$ . Compute:

$$\mathcal{Q}[\check{S}_1](x_3|x_2) = \mathcal{Q}[A_{1,1}](x_3|x_2) = \sum_{x_1} p(x_3|x_1, x_2) \cdot p(x_1).$$

d)  $S_1 = \{v_3\} = \check{S}_1$ . Compute:

$$\mathcal{Q}[S_1](x_3|x_2) = \mathcal{Q}[\check{S}_1](x_3|x_2) = \sum_{x_1} p(x_3|x_1, x_2) \cdot p(x_1).$$

5. For  $S_2 = \{v_2\}$ :

a)  $D_{2,0} = \text{Dist}^{[V]}(S_2) = \{v_2\}$ , see Figure 12 (a), (d). Compute:

$$\mathcal{Q}[D_{2,0}](x_2|x_1) = q(x_2|x_1) = \frac{\mathcal{Q}[V](x_1, x_2)}{\mathcal{Q}[V](x_1)} = p(x_2|x_1).$$

b)  $A_{2,1} = \text{Anc}^{[D_{2,0}]}(S_1) = \{v_2\} = D_{2,0}$ . thus  $\check{S}_2 = D_{2,0}$ . Compute:

$$\mathcal{Q}[\check{S}_2](x_2|x_1) = \mathcal{Q}[D_{2,0}](x_2|x_1) = p(x_2|x_1).$$

c)  $S_2 = \{v_2\} = \check{S}_2$ . Compute:

$$\mathcal{Q}[S_2](x_2|x_1) = \mathcal{Q}[\check{S}_2](x_2|x_1) = p(x_2|x_1).$$

6. Since both  $\check{S}_1 = S_1 = \{v_3\}$  and  $\check{S}_2 = S_2 = \{v_2\}$  we can compute:

$$\begin{aligned} \mathcal{Q}[B^C](x_2, x_3|x_1) &= \mathcal{Q}[S_2](x_2|x_1) \cdot \mathcal{Q}[S_1](x_3|x_2) \\ &= p(x_2|x_1) \cdot \sum_{x'_1} p(x_3|x'_1, x_2) \cdot p(x'_1), \\ p(x_3|\text{do}(x_1)) &= \sum_{x_2} \mathcal{Q}[B^C](x_2, x_3|x_1) \\ &= \sum_{x_2} p(x_2|x_1) \cdot \sum_{x'_1} p(x_3|x'_1, x_2) \cdot p(x'_1) \\ &= \sum_{x'_1, x'_2} p(x_3|x'_1, x'_2) \cdot p(x'_1) \cdot p(x'_2|x_1). \end{aligned}$$

So we can identify  $P(X_3|\text{do}(X_1))$  from  $P(X_1, X_2, X_3)$  via the mass function from above.

**Proofs - Soundness Criteria** We have seen in Corollary 5.3.18 that the *ID-Algorithm 5.3.17* is *sound up to oracle choices*. In this subsection we want to investigate the possibility of other forms of soundness that would allow for stronger forms of identifiability and/or trackability.

**Theorem 5.3.31** (Soundness up to null-sets). *Let  $G = (J, V, E, L)$  be a CADMG with a fixed topological order  $<$ . Consider the class of L-CBNs  $M$  with observable CADMG  $G$  such that the following holds:*

1. *The measurable spaces  $\mathcal{X}_v$  come equipped with a measure  $\mu_v$ ,  $v \in V$ ,*
2. *for every subset  $D \subseteq V$  the interventional Markov kernel  $\mathcal{Q}[D] = P(X_D|\text{do}(X_{J \cup V \setminus D}))$  is absolute continuous w.r.t. the product measure  $\mu_D := \otimes_{v \in D} \mu_v$  and vice versa:*

$$\mu_D \ll \mathcal{Q}[D] \ll \mu_D.$$

*If the ID-algorithm does not produce FAIL for input  $B, C \subseteq G$ , then  $P(X_B|\text{do}(X_{J \cup C}))$  is “almost-surely” trackable from  $P(X_V|\text{do}(X_J))$  and  $G$  for such CBNs  $M$ , i.e. the Markov kernel that was output by the ID-algorithm equals  $P(X_B|\text{do}(X_{J \cup C}))$  up to a  $\mu_{V \setminus B^C}$ -null set in  $\mathcal{X}_{J \cup V \setminus B^C}$ .*

*Proof.* Since by the second assumption we have for each  $\mu_v$  and any fixed value  $x_{\{v\}^c}$ :

$$P(X_v | \text{do}(X_{\{v\}^c} = x_{\{v\}^c})) \ll \mu_v \ll P(X_v | \text{do}(X_{\{v\}^c} = x_{\{v\}^c})),$$

we can w.l.o.g. assume that the  $\mu_v$  are probability measures for  $v \in V$ .

a) Now consider any subset  $A \subseteq V$  and  $D \in \mathcal{D}[A]$  and  $v \in D$ . We abbreviate:

$$G' := G_{\text{do}(V \setminus A)}, \quad A_{<} := \text{Pred}_{<}^{[A]}(v), \quad D_{<} := \text{Pred}_{<}^{[D]}(v).$$

Assume that we have a Markov kernel:

$$\mathcal{K}[A] = \tilde{P}(X_A | \text{do}(X_{J \cup V \setminus A})) : \mathcal{X}_{J \cup V \setminus A} \dashrightarrow \mathcal{X}_A,$$

such that:

$$\mathcal{K}[A] = \mathcal{Q}[A] \quad \mu_{V \setminus A}\text{-a.s.}$$

Note that the almost sure equality from above implies the equality:

$$\mathcal{K}[A] \otimes \mu_{V \setminus A} = \mathcal{Q}[A] \otimes \mu_{V \setminus A}. \quad (34)$$

We want to show that if we perform the steps of the ID-algorithm that computes  $\mathcal{Q}[D]$  from  $\mathcal{Q}[A]$  on  $\mathcal{K}[A]$  then the corresponding output, abbreviated as  $\mathcal{K}[D]$ , satisfies:

$$\mathcal{K}[D] = \mathcal{Q}[D] \quad \mu_{V \setminus D}\text{-a.s.}$$

For this consider any version of the conditional of the following marginal of  $\mathcal{K}[A]$ :

$$\tilde{P}(X_{A_{\leq}} | \text{do}(X_{J \cup V \setminus A})) \quad \text{w.r.t.} \quad \tilde{P}(X_{A_{<}} | \text{do}(X_{J \cup V \setminus A})),$$

which we will denote by:

$$K(X_v | X_{J \cup V \setminus A_{\geq}}).$$

This by definition will satisfy:

$$\tilde{P}(X_{A_{\leq}} | \text{do}(X_{J \cup V \setminus A})) = K(X_v | X_{J \cup V \setminus A_{\geq}}) \otimes \tilde{P}(X_{A_{<}} | \text{do}(X_{J \cup V \setminus A})). \quad (35)$$

This then implies:

$$\begin{aligned} & P(X_{A_{\leq}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\ & \stackrel{\text{Eq. 34}}{=} \tilde{P}(X_{A_{\leq}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\ & \stackrel{\text{Eq. 35}}{=} K(X_v | X_{J \cup V \setminus A_{\geq}}) \otimes \tilde{P}(X_{A_{<}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\ & \stackrel{\text{Eq. 34}}{=} K(X_v | X_{J \cup V \setminus A_{\geq}}) \otimes P(X_{A_{<}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}). \end{aligned}$$

This shows that  $K(X_v | X_{J \cup V \setminus A_{\geq}})$  is a version of the conditional of:

$$P(X_{A_{\leq}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \quad \text{w.r.t.} \quad P(X_{A_{<}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}).$$

Note that by the interventional ordered local Markov property, Proposition 5.3.10, there exists a Markov kernel  $Q(X_v|X_{\text{Mb}_{\xi'}(v)})$  that simultaneously is a version of both:

$$P(X_v|X_{A_{<}}, \text{do}(X_{J \cup V \setminus A})) \quad \text{and} \quad P(X_v|X_{D_{<}}, \text{do}(X_{J \cup V \setminus D})),$$

and is thus, in particular, another version of the conditional of:

$$P(X_{A_{\leq}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \quad \text{w.r.t.} \quad P(X_{A_{<}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}).$$

Now consider the (measurable) set where these two conditional Markov kernels deviate:

$$\tilde{N} := \left\{ x_{J \cup V \setminus A_{\geq}} \in \mathcal{X}_{J \cup V \setminus A_{\geq}} \mid K(X_v|X_{J \cup V \setminus A_{\geq}} = x_{J \cup V \setminus A_{\geq}}) \neq Q(X_v|X_{\text{Mb}_{\xi'}(v)} = x_{\text{Mb}_{\xi'}(v)}) \right\},$$

and  $N := \tilde{N} \times \mathcal{X}_{A_{\geq}} \subseteq \mathcal{X}_{J \cup V}$ . Since conditional Markov kernels are essentially unique we get that for every  $x_J \in \mathcal{X}_J$  we have:

$$(\mathcal{Q}[A] \otimes \mu_{V \setminus A})(N^{x_J}|x_J) = 0.$$

Since, by assumption, we have:  $\mu_A \ll \mathcal{Q}[A]$ , we get for every  $x_J \in \mathcal{X}_J$ :

$$(\mu_D \otimes \mu_{V \setminus D})(N^{x_J}) = \mu_V(N^{x_J}) = (\mu_A \otimes \mu_{V \setminus A})(N^{x_J}) = 0.$$

Since, by assumption, we also have:  $\mathcal{Q}[D] \ll \mu_D$ , we get for every  $x_J \in \mathcal{X}_J$ :

$$(\mathcal{Q}[D] \otimes \mu_{V \setminus D})(N^{x_J}|x_J) = 0.$$

Let  $\hat{N} := \tilde{N} \times \mathcal{X}_{A_{\geq} \setminus D_{\geq}}$ . Since  $D_{\geq} \subseteq A_{\geq}$  the set  $N$  is of the form:

$$N = \tilde{N} \times \mathcal{X}_{A_{\geq}} = \hat{N} \times \mathcal{X}_{D_{\geq}}.$$

So the above shows that we have for every  $x_J \in \mathcal{X}_J$ :

$$(P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D})) (\hat{N}^{x_J}|x_J) = (\mathcal{Q}[D] \otimes \mu_{V \setminus D})(N^{x_J}|x_J) = 0.$$

This shows that  $K(X_v|X_{J \cup V \setminus A_{\geq}})$  and  $Q(X_v|X_{\text{Mb}_{\xi'}(v)})$  agree up to a measurable  $(P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}))$ -null set. Remember that  $Q(X_v|X_{\text{Mb}_{\xi'}(v)})$  satisfies:

$$P(X_{D_{\leq}} | \text{do}(X_{J \cup V \setminus D})) = Q(X_v|X_{\text{Mb}_{\xi'}(v)}) \otimes P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})).$$

Together with the above we thus get:

$$\begin{aligned} & P(X_{D_{\leq}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\ &= Q(X_v|X_{\text{Mb}_{\xi'}(v)}) \otimes P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\ &= K(X_v|X_{J \cup V \setminus A_{\geq}}) \otimes P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}). \end{aligned}$$

This shows that  $K(X_v|X_{J \cup V \setminus A_{\geq}})$  is version of the conditional of:

$$P(X_{D_{\leq}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \quad \text{w.r.t.} \quad P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}).$$

If we let  $v$  run through  $D = \{v_1, \dots, v_K\}$ ,  $v_1 < v_2 < \dots < v_K$  in reverse topological order we inductively get:

$$\begin{aligned}
& \mathcal{Q}[D] \otimes \mu_{V \setminus D} \\
&= P(X_D | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\
&= K(X_{v_K} | X_{J \cup V \setminus A_{\geq v_K}}) \otimes P(X_{D_{< v_K}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\
&= K(X_{v_K} | X_{J \cup V \setminus A_{\geq v_K}}) \otimes P(X_{D_{\leq v_{K-1}}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\
&= K(X_{v_K} | X_{J \cup V \setminus A_{\geq v_K}}) \otimes K(X_{v_{K-1}} | X_{J \cup V \setminus A_{\geq v_{K-1}}}) \otimes P(X_{D_{< v_{K-1}}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\
&= \dots \\
&= \left( \bigotimes_{v \in D}^{\succ} K(X_v | X_{J \cup V \setminus A_{\geq v}}) \right) \otimes \mu_{V \setminus D}(X_{V \setminus D}).
\end{aligned}$$

Since such factorizations are essentially unique we get that:

$$\mathcal{Q}[D] = \bigotimes_{v \in D}^{\succ} K(X_v | X_{J \cup V \setminus A_{\geq v}}) =: \mathcal{K}[D] \quad \mu_{V \setminus D}\text{-a.s.}$$

This shows the claim.

b) We now reverse the situation. For a subset  $A \subseteq V$  and every  $D \in \mathcal{D}[A] = \{D_1, \dots, D_L\}$ , consider that we are given a Markov kernel:

$$\mathcal{K}[D] = \tilde{P}(X_D | \text{do}(X_{J \cup V \setminus D})) : \mathcal{X}_{J \cup V \setminus D} \dashrightarrow \mathcal{X}_D,$$

such that:

$$\mathcal{K}[D] = \mathcal{Q}[D] \quad \mu_{V \setminus D}\text{-a.s.},$$

which implies the equality:

$$\mathcal{K}[D] \otimes \mu_{V \setminus D} = \mathcal{Q}[D] \otimes \mu_{V \setminus D}. \quad (36)$$

We want to show that we then also have:

$$\mathcal{K}[A] := \left[ \bigotimes_{D \in \mathcal{D}[A]}^{\succ} \right] \mathcal{K}[D] = \mathcal{Q}[A] \quad \mu_{V \setminus A}\text{-a.s.}$$

For this fix a node  $v \in D$  and note that  $Q(X_v | X_{\text{Mb}_{\prec}^{\mathcal{G}'(v)}})$  satisfies:

$$P(X_{D_{\leq}} | \text{do}(X_{J \cup V \setminus D})) = Q(X_v | X_{\text{Mb}_{\prec}^{\mathcal{G}'(v)}}) \otimes P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})). \quad (37)$$

We then get the equalities:

$$\begin{aligned}
& \tilde{P}(X_{D_{\leq}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\
& \stackrel{\text{Eq. 36}}{=} P(X_{D_{\leq}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\
& \stackrel{\text{Eq. 37}}{=} Q(X_v | X_{\text{Mb}_{\prec}^{\mathcal{G}'(v)}}) \otimes P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}) \\
& \stackrel{\text{Eq. 36}}{=} Q(X_v | X_{\text{Mb}_{\prec}^{\mathcal{G}'(v)}}) \otimes \tilde{P}(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D})) \otimes \mu_{V \setminus D}(X_{V \setminus D}).
\end{aligned}$$

So  $Q(X_v|X_{\text{Mb}_{<}^{\mathcal{G}'(v)}})$  is a conditional of:

$$\tilde{P}(X_{D_{\leq}}|\text{do}(X_{J\cup V\setminus D})) \otimes \mu_{V\setminus D}(X_{V\setminus D}) \quad \text{w.r.t.} \quad \tilde{P}(X_{D_{<}}|\text{do}(X_{J\cup V\setminus D})) \otimes \mu_{V\setminus D}(X_{V\setminus D}),$$

that does not depend on  $X_{A_{\geq}}$  (as  $\text{Mb}_{<}^{\mathcal{G}'(v)} \subseteq A_{<}$ ).

Now consider any other version of the conditional of:

$$\tilde{P}(X_{D_{\leq}}|\text{do}(X_{J\cup V\setminus D})) \otimes \mu_{V\setminus D}(X_{V\setminus D}) \quad \text{w.r.t.} \quad \tilde{P}(X_{D_{<}}|\text{do}(X_{J\cup V\setminus D})) \otimes \mu_{V\setminus D}(X_{V\setminus D}),$$

that does not depend on variables attached to  $A_{\geq}$  and which we will denote by:

$$K(X_v|X_{J\cup V\setminus A_{\geq}}).$$

Note that such a Markov kernel exists, as  $Q(X_v|X_{\text{Mb}_{<}^{\mathcal{G}'(v)}})$  is such one.

The same argumentation with  $K(X_v|X_{J\cup V\setminus A_{\geq}})$  in place of  $Q(X_v|X_{\text{Mb}_{<}^{\mathcal{G}'(v)}})$ , using Eq. 36, shows that both,  $Q(X_v|X_{\text{Mb}_{<}^{\mathcal{G}'(v)}})$  and  $K(X_v|X_{J\cup V\setminus A_{\geq}})$ , are then conditionals of

$$P(X_{D_{\leq}}|\text{do}(X_{J\cup V\setminus D})) \otimes \mu_{V\setminus D}(X_{V\setminus D}) \quad \text{w.r.t.} \quad P(X_{D_{<}}|\text{do}(X_{J\cup V\setminus D})) \otimes \mu_{V\setminus D}(X_{V\setminus D}),$$

that do not depend on  $X_{A_{\geq}}$ . Now let:

$$\tilde{N} := \left\{ x_{J\cup V\setminus A_{\geq}} \in \mathcal{X}_{J\cup V\setminus A_{\geq}} \mid K(X_v|X_{J\cup V\setminus A_{\geq}} = x_{J\cup V\setminus A_{\geq}}) \neq Q(X_v|X_{\text{Mb}_{<}^{\mathcal{G}'(v)}} = x_{\text{Mb}_{<}^{\mathcal{G}'(v)}}) \right\},$$

and  $N := \tilde{N} \times \mathcal{X}_{A_{\geq}} \subseteq \mathcal{X}_{J\cup V}$ . Again, since both are versions of the same conditional we get:

$$(\mathcal{Q}[D] \otimes \mu_{V\setminus D})(N^{x_J}|x_J) = 0,$$

for every  $x_J \in \mathcal{X}_J$ . Since  $\mathcal{Q}[D] \ll \mu_D$  we get for every  $x_J \in \mathcal{X}_J$ :

$$(\mu_A \otimes \mu_{V\setminus A})(N^{x_J}) = \mu_V(N^{x_J}) = (\mu_D \otimes \mu_{V\setminus D})(N^{x_J}) = 0.$$

Since also  $\mathcal{Q}[A] \ll \mu_A$  we get for every  $x_J \in \mathcal{X}_J$ :

$$(\mathcal{Q}[A] \otimes \mu_{V\setminus A})(N^{x_J}|x_J) = 0.$$

Since  $N = \tilde{N} \times \mathcal{X}_{A_{\geq}}$  we get for every  $x_J \in \mathcal{X}_J$ :

$$(P(X_{A_{<}}|\text{do}(X_{J\cup V\setminus A})) \otimes \mu_{V\setminus A}(X_{V\setminus A}))(N^{x_J}|x_J) = (\mathcal{Q}[A] \otimes \mu_{V\setminus A})(N^{x_J}|x_J) = 0.$$

This shows that the Markov kernels  $Q(X_v|X_{\text{Mb}_{<}^{\mathcal{G}'(v)}})$  and  $K(X_v|X_{J\cup V\setminus A_{\geq}})$  are equal up to some  $(P(X_{A_{<}}|\text{do}(X_{J\cup V\setminus A})) \otimes \mu_{V\setminus A}(X_{V\setminus A}))$ -null set. Note that with this we get the



factorization, using  $A = \{v_1, \dots, v_K\}$ ,  $v_1 < \dots < v_K$ :

$$\begin{aligned}
& \mathcal{Q}[A] \otimes \mu_{V \setminus A} \\
&= P(X_{A \leq v_K} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\
&= Q(X_{v_K} | X_{\text{Mb}_{\xi'}(v_K)}) \otimes P(X_{A < v_K} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\
&= K(X_{v_K} | X_{J \cup V \setminus A \geq v_K}) \otimes P(X_{A < v_K} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\
&= K(X_{v_K} | X_{J \cup V \setminus A \geq v_K}) \otimes P(X_{A \leq v_{K-1}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\
&= K(X_{v_K} | X_{J \cup V \setminus A \geq v_K}) \otimes Q(X_{v_K} | X_{\text{Mb}_{\xi'}(v_K)}) \otimes P(X_{A < v_{K-1}} | \text{do}(X_{J \cup V \setminus A})) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\
&= \dots \\
&= \left( \bigotimes_{v \in A}^{\succ} K(X_v | X_{J \cup V \setminus A \geq v}) \right) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\
&= \left( \left[ \bigotimes_{D \in \mathcal{D}[A]}^{\succ} \right] \mathcal{K}[D] \right) \otimes \mu_{V \setminus A}(X_{V \setminus A}) \\
&= \mathcal{K}[A] \otimes \mu_{V \setminus A}(X_{V \setminus A}).
\end{aligned}$$

Since such factorizations are essentially unique we get:

$$\mathcal{K}[A] = \mathcal{Q}[A] \quad \mu_{V \setminus A}\text{-a.s.}$$

This shows the claim.

c) Now let  $D \subseteq V$  and  $A \subseteq D$  with  $A = \text{Anc}^{[D]}(A)$ . Consider that we are given a Markov kernel:

$$\mathcal{K}[D] = \tilde{P}(X_D | \text{do}(X_{J \cup V \setminus D})) : \mathcal{X}_{J \cup V \setminus D} \dashrightarrow \mathcal{X}_D,$$

such that:

$$\mathcal{K}[D] = \mathcal{Q}[D] \quad \mu_{V \setminus D}\text{-a.s.},$$

which implies the equality:

$$\mathcal{K}[D] \otimes \mu_{V \setminus D} = \mathcal{Q}[D] \otimes \mu_{V \setminus D}.$$

We want to show that the  $A$ -marginal of  $\mathcal{K}[D]$  equals  $\mathcal{Q}[A]$  up to  $\mu_{V \setminus A}$ -null set.

For this let  $\mathcal{K}[A]$  be the  $A$ -marginal of  $\mathcal{K}[D]$ :

$$\mathcal{K}[A] := \tilde{P}(X_A | \text{do}(X_{J \cup V \setminus D})) : \mathcal{X}_{J \cup V \setminus A} \dashrightarrow \mathcal{X}_A.$$

Note that  $\mathcal{Q}[A]$  is the  $A$ -marginal of  $\mathcal{Q}[D]$ . Marginalizing out  $X_{D \setminus A}$  on both sides in the above equation gives us:

$$\mathcal{K}[A] \otimes \mu_{V \setminus D} = \mathcal{Q}[A] \otimes \mu_{V \setminus D}.$$

Multiplying both sides with  $\mu_{D \setminus A}$  gives:

$$\mathcal{K}[A] \otimes \mu_{V \setminus A} = \mathcal{K}[A] \otimes \mu_{V \setminus D} \otimes \mu_{D \setminus A} = \mathcal{Q}[A] \otimes \mu_{V \setminus D} \otimes \mu_{D \setminus A} = \mathcal{Q}[A] \otimes \mu_{V \setminus A}.$$

Since such factorizations are essentially unique we get:

$$\mathcal{K}[A] = \mathcal{Q}[A] \quad \mu_{V \setminus A}\text{-a.s.}$$

This shows the claim.

This covers all cases of the ID-algorithm and thus shows the claim.  $\square$

**Theorem 5.3.32** (Soundness up to continuous choices for strictly positive CBNs). *Let  $G = (J, V, E, L)$  be a CADMG with a fixed topological order  $<$ . Consider the class of L-CBNs  $M$  with observable CADMG  $G$  such that the following holds:*

1. *The spaces  $\mathcal{X}_v$  are Polish spaces for  $v \in J \cup V$ ,*
2. *for every subset  $D \subseteq V$  the interventional Markov kernel  $\mathcal{Q}[D] = P(X_D | \text{do}(X_{J \cup V \setminus D}))$  is strictly positive (on non-empty open subsets of  $\mathcal{X}_D$ ), and:*
3. *for every  $v \in D$  the Markov kernel  $Q(X_v | X_{\text{Mb}_{<}^{G, \text{do}(D^c)}(v)})$  can be chosen to be continuous, viewed as a map:  $\mathcal{X}_{\text{Mb}_{<}^{G, \text{do}(D^c)}(v)} \rightarrow \mathcal{P}(\mathcal{X}_v)$ .*

*If the ID-algorithm does not produce FAIL for input  $B, C \subseteq G$ , then  $P(X_B | \text{do}(X_{J \cup C}))$  is identifiable and trackable “up to continuous choices of conditional Markov kernels” from  $P(X_V | \text{do}(X_J))$  and  $G$  for such CBNs  $M$ , i.e. if every occurring conditional Markov kernel is chosen to be continuous (which will always be possible by the assumptions made).*

*Proof.* For a district  $D \in \mathcal{D}[V]$  and  $v \in D$ , by assumption, there exists a *continuous* version of  $Q(X_v | X_{\text{Mb}_{<}^G(v)})$ , which is also a version of:

$$P(X_v | X_{\text{Pred}_{<}^{[V]}(v)}, \text{do}(X_J)) \quad \text{and} \quad P(X_v | X_{\text{Pred}_{<}^{[D]}(v)}, \text{do}(X_{J \cup V \setminus D})).$$

We abbreviate  $V_{<} := \text{Pred}_{<}^{[V]}(v)$  and  $D_{<} := \text{Pred}_{<}^{[D]}(v)$  in the following.

Now consider any *continuous* version of the conditional Markov kernel  $P(X_v | X_{V_{<}}, \text{do}(X_J))$ . Note that such a version always exists because  $Q(X_v | X_{\text{Mb}_{<}^G(v)})$  is already an existing continuous version. Since  $P(X_{V_{<}} | \text{do}(X_J))$  is strictly positive, as the marginal of  $\mathcal{Q}[V]$ , Lemma 2.4.23 implies then the “sure” equality:

$$P(X_v | X_{V_{<}}, \text{do}(X_J)) = Q(X_v | X_{\text{Mb}_{<}^G(v)}).$$

So any continuous version of  $P(X_v | X_{V_{<}}, \text{do}(X_J))$  necessarily agrees with  $Q(X_v | X_{\text{Mb}_{<}^G(v)})$  on all points.

Similarly, using the same arguments, we get that  $Q(X_v | X_{\text{Mb}_{<}^G(v)})$  is “surely” equal to every *continuous* version of the conditional  $P(X_v | X_{D_{<}}, \text{do}(X_{J \cup V \setminus D}))$ , as also the Markov kernel  $P(X_{D_{<}} | \text{do}(X_{J \cup V \setminus D}))$  is strictly positive, as a marginal of  $\mathcal{Q}[D]$ . So we get:

$$Q(X_v | X_{\text{Mb}_{<}^G(v)}) = P(X_v | X_{D_{<}}, \text{do}(X_{J \cup V \setminus D})).$$

This means that if we pick a/the *continuous* version of the conditional  $P(X_v|X_{V<}, \text{do}(X_J))$  then it is “surely” equal to the/every *continuous* version of the conditional  $P(X_v|X_{D<}, \text{do}(X_{J \cup V \setminus D}))$ :

$$P(X_v|X_{V<}, \text{do}(X_J)) = Q(X_v|X_{\text{Mb}_{\leq}(v)}^G) = P(X_v|X_{D<}, \text{do}(X_{J \cup V \setminus D})).$$

These arguments, repeated for subgraphs, then show that in the ID-algorithm for every occuring conditional and product (e.g. for districts and the final product) we end up with distinct and correct choices for all Markov kernels. This then also shows the identifiability of such CBNs  $M$  (in the not-FAIL case).  $\square$

**Lemma 5.3.33.** *Consider an L-CBN:*

$$M = \left( G^+ = (J, (V, U), E^+), \left( P_v(X_v | \text{do}(X_{\text{Pa}^{G^+}(v)})) \right)_{v \in V \cup U} \right),$$

with observable CADMG  $G = (J, V, E, L)$  and fixed topological order  $<$ . Assume that for every  $v \in V$  we have a measure  $\mu_v$  on  $\mathcal{X}_v$  such that  $P_v(X_v | \text{do}(X_{\text{Pa}^{G^+}(v)}))$  has a (strictly positive) density w.r.t.  $\mu_v$ :

$$p(x_v | \text{do}(x_{\text{Pa}^{G^+}(v)})) > 0.$$

Furthermore, we put for  $x_V \in \mathcal{X}_V$ ,  $x_U \in \mathcal{X}_U$ ,  $x_J \in \mathcal{X}_J$ :

$$p(x_V | x_U, \text{do}(x_J)) := \prod_{v \in V} p(x_v | \text{do}(x_{\text{Pa}^{G^+}(v)})),$$

and then integrate in reverse order of  $<$ :

$$p(x_V | \text{do}(x_J)) := \int \cdots \int_{\mathcal{X}_U} p(x_V | x_U, \text{do}(x_J)) \bigotimes_{u \in U}^> P_u(X_u \in dx_u | \text{do}(X_{\text{Pa}^{G^+}(v)} = x_{\text{Pa}^{G^+}(v)})).$$

Then the former is a (strictly positive) density of  $P(X_V, X_U | \text{do}(X_J))$  w.r.t.

$$\bigotimes_{v \in V}^> \mu_v \otimes \bigotimes_{u \in U}^> P_u(X_u | \text{do}(X_{\text{Pa}^{G^+}(v)})),$$

and the latter a (strictly positive) density of  $P(X_V | \text{do}(X_J))$  w.r.t.

$$\mu_V := \bigotimes_{v \in V} \mu_v.$$

Similarly, for every  $D \subseteq V$  the interventional Markov kernel  $P(X_D | \text{do}(X_{J \cup V \setminus D}))$  has a (strictly positive) density w.r.t.  $\mu_D := \bigotimes_{v \in D} \mu_v$ .

*Proof.* The claim can be shown by integrating the above densities over product sets  $A = \prod_v A_v$ . Inductively we can use Fubini’s theorem and:

$$\int_{A_v} p(x_v | \text{do}(x_{\text{Pa}^{G^+}(v)})) \mu_v(dx_v) = P_v(X_v \in A_v | \text{do}(X_{\text{Pa}^{G^+}(v)} = x_{\text{Pa}^{G^+}(v)})).$$

Regarding strict positivity, note that if  $f(x) > 0$  for all  $x$ , then  $\int f d\mu > 0$  for non-trivial  $\mu$ . So strict positivity is preserved through integration.  $\square$

## 6. Structural Causal Models

Structural Causal Models (SCMs), also known as Structural Equation Models (SEMs), or Non-Parametric Structural Equation Models (NP-SEMs), provide a class of causal models that can model causal cycles. SCMs trace back to the early work on path analysis by geneticist Sewall Wright [Wri21], made their way to econometrics [Haa43, SW60], and became popular in AI due to the work of Judea Pearl [Pea09] and many others. In these lecture notes, we give a modern treatment inspired by our own research on the matter [BFPM21, FM20].

### 6.1. Motivation

While causal Bayesian networks (with input nodes and latent variables) provide a powerful causal modeling class, there is an important aspect of causality that cannot be modeled with causal Bayesian networks, namely causal *cycles*. For example, increasing temperature at the poles may cause sea ice to melt, which leads to more absorption of sunlight because white ice is replaced by blue sea water, which in turn leads to further temperature increase (see also Figure 13(a)). Because a causal Bayesian network is acyclic by definition, such a model can only be described by a causal Bayesian network by introducing multiple variables corresponding with measurements of the same quantities at different points in time (Figure 13(b)). In contrast, an SCM can directly represent causal cycles and is often appropriate for modeling systems with feedback loops that are stable, i.e., where negative feedback dominates potential positive feedback. An illustrative example is a system composed of different masses connected via springs in an environment with friction (see also Section 6.10).

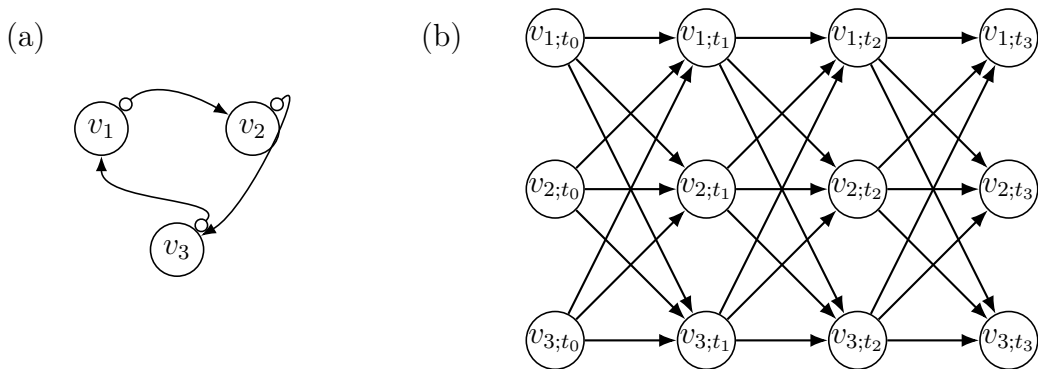


Figure 13: (a) Directed Graph (DG) representing a causal cycle. As an example,  $v_1$  could be the average temperature in a certain area at the North pole,  $v_2$  the amount of sea ice present in the area, and  $v_3$  the amount of sunlight absorbed in the area. This gives an example of a positive (self-reinforcing) feedback loop. (b) Alternative Directed Acyclic Graph (DAG) where the variables correspond with the same quantities but measured at different time points  $t_0 < t_1 < t_2 < t_3$ .

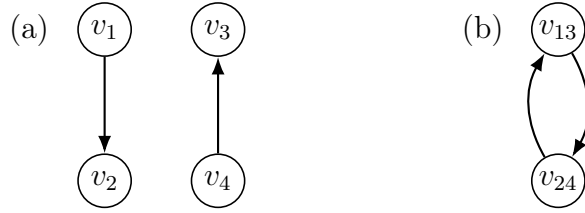


Figure 14: (a) ADMG with output nodes  $v_1, v_2, v_3, v_4$  corresponding with endogenous variables  $X_1, X_2, X_3, X_4$ . (b) DMG corresponding to a coarser representation obtained by merging variables into  $X_{13} := (X_1, X_3)$  and  $X_{24} := (X_2, X_4)$ .

Even if a ‘fine-grained’ causal model is acyclic, merging variables may introduce cycles at a more ‘coarse-grained’ level of description.

**Example 6.1.1.** *Suppose that we have a CBN with the ADMG in Figure 14(a), representing four variables  $X_1, X_2, X_3, X_4$ . If we chose an alternative representation in terms of pairs  $X_{13} := (X_1, X_3)$  and  $X_{24} := (X_2, X_4)$ , then we would end up with a CBN with the DMG in Figure 14(b). However, that is a contradiction as the graph of a CBN is acyclic by definition.*

This example shows that the class of CBNs is not closed under the operation of merging variables. The class of (simple) SCMs to be introduced later is actually closed under the operation of merging variables.

Finally, there exist systems in which the directionality of causal relations is context-dependent.

**Example 6.1.2.** *Consider Ohm’s law  $V = IR$  (voltage equals current times resistance) to model the voltage across and current through a resistance. If we connect the resistance to a voltage source, the voltage determines the current. If we connect the resistance to a current source, then it is the other way around: the current determines the voltage. Both cases separately can be modeled with a CBN (Figure 15(a–b), respectively). If we let a coin flip determine which of the two sources the resistance is connected to, we obtain a mixture which cannot be modeled as a single CBN (Figure 15(c)).*

Similar behavior is often encountered in complex systems in biology, chemistry, engineering and economy. This is yet another motivation to extend the causal modeling framework to allow for cycles.

In this chapter, we will introduce the class of SCMs, which generalize CBNs to allow for cycles, which allows us to deal elegantly with all motivating examples discussed here.

## 6.2. Solving SCMs

An SCM is specified in terms of (deterministic) functions and distributions, rather than in terms of Markov kernels (stochastic functions).

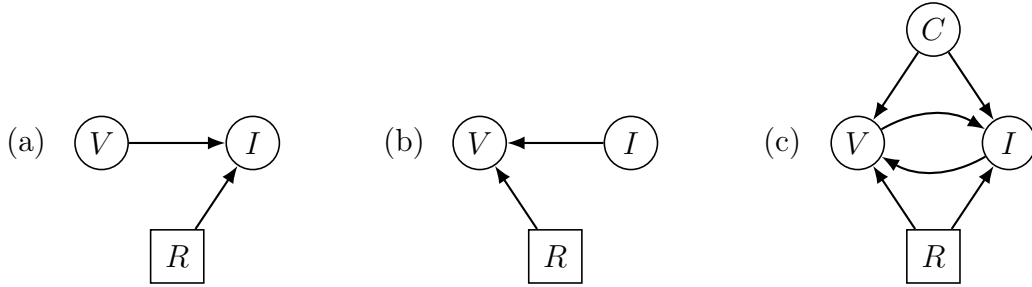


Figure 15: Different causal models corresponding to modeling the current through a resistance using Ohm’s law. (a) Voltage causes current. (b) Current causes voltage. (c) Mixture model where the causal relationship between voltage and current depends on the result of a coin flip.

In contrast with many definitions encountered in the literature,<sup>28</sup> we will explicitly distinguish three types of variables: exogenous random variables, exogenous input variables, and endogenous variables.

**Definition 6.2.1** (Structural Causal Model). *A Structural Causal Model is a tuple  $M = (J, V, W, \mathcal{X}, P, f)$  such that*

- $J, V, W$  are disjoint finite sets of labels for the exogenous input variables, the endogenous variables and the exogenous random variables, respectively;
- the domain  $\mathcal{X} = \prod_{i \in J \cup V \cup W} \mathcal{X}_i$  is a product of standard measurable spaces  $\mathcal{X}_i$ ;
- the exogenous distribution  $P$  is a probability distribution on  $\mathcal{X}_W$  that factorizes as a product  $P = \bigotimes_{w \in W} P_w$  of probability distributions  $P_w \in \mathcal{P}(\mathcal{X}_w)$ ;
- the causal mechanism is specified by the measurable function  $f : \mathcal{X} \rightarrow \mathcal{X}_V$ .

Often, the causal mechanism  $f$  and the exogenous distribution  $P$  depend (in a measurable way) on exogenous parameters  $\theta \in \Theta$ , which we may make explicit by writing  $f_\theta$  and  $P_\theta$  instead, giving a parameterized SCM  $M_\theta = (J, V, W, \mathcal{X}, P_\theta, f_\theta)$ . The family  $(M_\theta)_{\theta \in \Theta}$  is then an SCM family.<sup>29</sup>

One can also think about an SCM as describing an input/output system, with free inputs  $J$ , random inputs  $W$  with distribution  $P$ , outputs  $V$  and input/output mechanism  $f$ . Structural causal models can be regarded as a marriage of statistical models as traditionally used in statistics (a parameterized family of distributions) with deterministic causal models (deterministic input/output systems) that are used informally in disciplines like physics and engineering.

<sup>28</sup>For example, [Pea09] only formally distinguishes exogenous random variables and endogenous variables.

<sup>29</sup>In line with the convention in machine learning, the word “model” refers to a SCM with a fixed choice of the parameters, and “model family” to a family of models indexed measurably by parameters. This contrasts with the terminology in statistics, where a family of distributions indexed measurably by parameters is called a “statistical model”.

**Remark 6.2.2.** *There are three crucial assumptions embodied in the modeling approach using SCMs:*

1. *The distinction between endogenous and exogenous variables: exogenous variables (i.e., exogenous input variables, exogenous random variables, and exogenous parameters) are **not caused** by endogenous variables, by assumption;*
2. *Exogenous random variables are **mutually independent**, and independent of the exogenous input variables in the sense that their probability distribution does not depend on the joint value of all exogenous input variables; however, we do allow for “dependencies” between exogenous input variables;*
3. *Exogenous parameters are distinguished from exogenous random variables in that the former describe “population” properties whereas the latter describe “individual” quantities.*

**Remark 6.2.3.** *Not all types of variables need to be present. Rather than giving separate definitions for ‘degenerate’ cases, we can stay in the formalism by defining what happens for empty label sets. For example, suppose the SCM is deterministic, i.e.,  $W = \emptyset$ . Then  $\mathcal{X}_W$  is an empty product (i.e., a product over 0 spaces), and by definition becomes a space  $\ast = \{\ast\}$  with a single element  $\ast$ , with the trivial sigma algebra  $\{\emptyset, \{\ast\}\}$ . The only possible probability distribution on such a space is the trivial distribution, i.e.,  $P(\{\ast\}) = 1$ . Similarly, it often happens that there are no exogenous input variables (if  $J$  is empty).*

**Remark 6.2.4.** *Often the exogenous random variables are latent and thought of as “noise”. However, since it may depend on the context whether a variable is observed or latent (e.g., in a training data set, the prediction target is typically observed, whereas it is latent in a test data set), we will not formally incorporate into the model any assumptions regarding which variables are observed and which are latent. In this aspect, our exposition deviates from many accounts on SCMs in the literature.<sup>30</sup>*

### 6.2.1. Potential Outcomes

An SCM defines *structural equations*, which are used to define the *potential outcomes* of the SCM.

**Definition 6.2.5** (Potential outcomes, structural equations). *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM and  $x_J \in \mathcal{X}_J$  an input value. A random variable  $X_{V \cup W}^{\text{do}(x_J)}$  with codomain  $\mathcal{X}_V \times \mathcal{X}_W$  is called a potential outcome of  $M$  for input  $x_J$  if the following two conditions hold:<sup>31</sup>*

1. *its  $W$ -component has the exogenous distribution specified by  $M$ :*

$$X_W^{\text{do}(x_J)} \sim P,$$

---

<sup>30</sup>For example, [Pea09] assumes exogenous (random) variables to be latent, and endogenous variables to be observed.

<sup>31</sup>Another notation for potential outcomes, commonly encountered in the literature, is  $X_{V \cup W}(x_J)$ .

2. it satisfies the structural equations entailed by  $M$  for input  $x_J$ :

$$X_V^{\text{do}(x_J)} = f(x_J, X_V^{\text{do}(x_J)}, X_W^{\text{do}(x_J)}) \text{ a.s..} \quad (38)$$

In case  $J = \emptyset$  we also write  $X_{V \cup W} := X_{V \cup W}^{\text{do}(\ast)}$  and refer to it simply as an outcome of  $M$ .

The SCM encodes the probability distributions of its (potential) outcomes. However, the structural equations (38) may have no solution or may have multiple different solutions. Therefore, even if a (potential) outcome exists (for a given input), it could be that its distribution is not uniquely determined by the SCM.

**Example 6.2.6.** Consider an SCM with parameters  $\alpha, \beta \in \mathbb{R}$ , endogenous real-valued variables  $X_1, X_2$ , real-valued exogenous input  $X_3$ , and structural equations

$$\begin{cases} X_1^{\text{do}(x_3)} = \alpha X_2^{\text{do}(x_3)} \\ X_2^{\text{do}(x_3)} = \beta X_1^{\text{do}(x_3)} + x_3. \end{cases}$$

If  $\alpha\beta = 1$ , then  $(X_1^{\text{do}(x_3=0)}, X_2^{\text{do}(x_3=0)}) = (\alpha x, x)$  is a potential outcome for input  $x_3 = 0$  for any value of  $x \in \mathbb{R}$ . Any mixture of these potential outcomes is also a potential outcome for input  $x_3 = 0$ . If  $\alpha\beta = 1$ , then for input  $x_3 \neq 0$ , the SCM admits no potential outcomes. If  $\alpha\beta \neq 1$ , then the potential outcomes are unique and given by  $(X_1^{\text{do}(x_3)}, X_2^{\text{do}(x_3)}) = (\frac{\alpha x_3}{1-\alpha\beta}, \frac{x_3}{1-\alpha\beta})$ .

In practice, an SCM is often specified more informally by writing down the corresponding structural equations and by giving the exogenous distribution of  $X_W$ . Any variables appearing on the r.h.s. of some structural equation that do not correspond with a structural equation for which that variable appears on the l.h.s., nor have a specified distribution, are then implicitly taken as exogenous inputs or parameters.

**Example 6.2.7** (Linear regression model with fixed design for the effect in terms of its cause). Assume that  $Y = \alpha X + \beta + \epsilon$  with  $Y \in \mathbb{R}$  representing the effect,  $X \in \mathbb{R}$  the cause,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  independent normally distributed measurement noise, and  $\alpha, \beta, \sigma^2 \in \mathbb{R}$  parameters.

This can be understood as the specification of an SCM  $M_{(\alpha, \beta, \sigma^2)} = (J, V, W, \mathcal{X}, P_{\sigma^2}, f_{\alpha, \beta})$  with  $J = \{X\}$ ,  $V = \{Y\}$ ,  $W = \{\epsilon\}$ ,  $\mathcal{X} = \mathbb{R}^3$ , exogenous distribution  $P_{\sigma^2}(X_\epsilon) = \mathcal{N}(0, \sigma^2)$  and causal mechanism  $f_{\alpha, \beta} : \mathbb{R}^3 \rightarrow \mathbb{R} : (x, y, \epsilon) \mapsto \alpha x + \beta + \epsilon$ .

If  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ , then  $(Y^{\text{do}(x)}, \epsilon^{\text{do}(x)}) := (\alpha x + \beta + \epsilon, \epsilon)$  is a potential outcome of  $M$ . If  $\epsilon$  is latent (which is usually the case), we just refer to  $Y^{\text{do}(x)}$  as the potential outcome.

## 6.2.2. Solutions

The language of conditional random variables allows us to give a neat definition of the *solution* of an SCM, which can also be thought of as a measurable family of (potential) outcomes with a shared underlying probability space.



**Definition 6.2.8** (Solutions of an SCM). Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. Let  $(\mathcal{U} \times \mathcal{X}_J, K(U|X_J))$  be a transition probability space and  $X : \mathcal{U} \times \mathcal{X}_J \rightarrow \mathcal{X}$  be a conditional random variable. If its push-forward Markov kernel  $K(X|X_J)$  satisfies

$$K(X_W|X_J) = P(X_W),$$

$$K(X_J|X_J) = \delta(X_J|X_J)$$

and  $X$  satisfies the structural equations

$$X_V = f(X_J, X_V, X_W) \quad K(X|X_J)\text{-a.s.} \quad (39)$$

then  $X$  is called a solution of  $M$ . Since the  $X_J$  component is trivial, we also refer to the component

$$X_{V \cup W} : \mathcal{U} \times \mathcal{X}_J \rightarrow \mathcal{X}_{V \cup W}$$

as a solution of  $M$ .

**Remark 6.2.9.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM and  $X : \mathcal{U} \times \mathcal{X}_J \rightarrow \mathcal{X}_{V \cup W}$  be a solution of  $M$ . Then for any  $x_J \in \mathcal{X}_J$ ,

$$X_{V \cup W}^{\text{do}(x_J)} : \mathcal{U} \rightarrow \mathcal{X}_{V \cup W} : u \mapsto X(u, x_J)$$

is a (potential) outcome of  $M$  for input  $x_J$ .

Not every SCM has solutions. Also, if they exist, solutions are not necessarily unique, even if they have the same underlying transition probability space.

**Example 6.2.10.** Consider an SCM with parameters  $\alpha, \beta, \mu, \sigma \in \mathbb{R}$ , endogenous real-valued variables  $X_1, X_2$ , exogenous random real-valued variable  $W_1$ , structural equations

$$\begin{cases} X_1 = \alpha X_2 \\ X_2 = \beta X_1 + W_1, \end{cases}$$

and exogenous distribution  $\mathcal{N}(\mu, \sigma^2)$ . If  $\alpha\beta = 1$ ,  $\mu = 0$  and  $\sigma^2 = 0$ , then  $(X_1, X_2, W_1) = (\alpha x, x, W_1)$  is a solution for any  $x \in \mathbb{R}$  and  $W_1 \sim \mathcal{N}(\mu, \sigma^2)$ . Any mixture of those solutions is also a solution in that case. If  $\alpha\beta = 1$  and  $\mu \neq 0$  or  $\sigma^2 \neq 0$ , then the SCM admits no solutions. If  $\alpha\beta \neq 1$ , then all solutions  $(X_1, X_2, W_1)$  satisfy  $(X_1, X_2) = (\frac{\alpha W_1}{1-\alpha\beta}, \frac{W_1}{1-\alpha\beta})$  a.s. and  $W_1 \sim \mathcal{N}(\mu, \sigma^2)$ .

The following remark relates the terminology to the cases most often considered in the literature.

**Remark 6.2.11.** If  $J = \emptyset$ , a solution of an SCM can be identified with a random variable  $X_{V \cup W}$  with codomain  $\mathcal{X}_V \times \mathcal{X}_W$  such that  $X_W \sim P$  and that satisfies the structural equations:

$$X_v = f_v(X_V, X_W) \quad \text{a.s.} \quad (40)$$

for each  $v \in V$ .

### 6.2.3. Markov Kernels of Solutions

Each solution of an SCM “has” a Markov kernel (similarly to how each random variable has a distribution).

**Notation 6.2.12.** *Let conditional random variable  $X : \mathcal{U} \times \mathcal{X}_J \rightarrow \mathcal{X}$  on transition space  $(\mathcal{U} \times \mathcal{X}_J, K(U|X_J))$  be a solution of an SCM  $M = (J, V, W, \mathcal{X}, P, f)$ . Its push-forward*

$$P(X \mid \text{do}(X_J)) := K(X|X_J) = X_*K(U|X_J)$$

*is a Markov kernel  $\mathcal{X}_J \dashrightarrow \mathcal{X}$  that we refer to as the Markov kernel of  $M$  corresponding to  $X$ . Since the  $J$ -component is trivial, we also refer to its marginal  $P(X_{V \cup W} \mid \text{do}(X_J))$  as such.*

Not all solutions of an SCM may yield the same Markov kernel (see also Example 6.2.10). Therefore, even if the SCM  $M$  is specified, the notation “ $P(X \mid \text{do}(X_J))$ ” is ambiguous, since it does not specify which solution the Markov kernel comes from. Because we will mostly restrict attention to so-called ‘simple’ SCMs for which the Markov kernel turns out to be unique, we will not worry about this.

**Remark 6.2.13.** *In case  $J = \emptyset$ , the Markov kernel  $P(X \mid \text{do}(X_J))$  corresponding to a solution  $X$  can be identified with its distribution  $P(X)$ , and one often refers to it as the distribution of  $M$  corresponding to  $X$ .*

### 6.2.4. Solution functions

We can construct solutions and (potential) outcomes in terms of solution functions of an SCM:

**Definition 6.2.14** (Solution function of an SCM). *Given an SCM  $M = (J, V, W, \mathcal{X}, P, f)$ , we call a measurable function  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  a solution function of  $M$  if for all  $x_W \in \mathcal{X}_W$ , for all  $x_J \in \mathcal{X}_J$ ,  $g(x_J, x_W)$  satisfies the structural equations.<sup>32</sup>*

$$g(x_J, x_W) = f(x_J, g(x_J, x_W), x_W).$$

**Remark 6.2.15.** *If  $g$  is a solution function for SCM  $M = (J, V, W, \mathcal{X}, P, f)$ , then for any random variable  $X_W$  with codomain  $\mathcal{X}_W$  and distribution  $P$ :*

1.  $X_{V,W}^{\text{do}(x_J)} := (g(x_J, X_W), X_W)$  is a (potential) outcome for  $M$  for input  $x_J \in \mathcal{X}_J$ ;
2.  $X_{V,W} := (g(X_J, X_W), X_W)$  is a solution of  $M$ ;
3. its corresponding Markov kernel is the push-forward

$$(g, \text{id}_{\mathcal{X}_W})_*(P \otimes \delta(X_J|X_J)).$$

---

<sup>32</sup>One can weaken this requirement by replacing the quantifiers by “for all  $x_J \in \mathcal{X}_J$ , for  $P$ -almost all  $x_W \in \mathcal{X}_W$ ” but we will not do so here, as the additional generality of the resulting theory comes at the cost of more technicalities in most definitions and proofs.

Not all (potential) outcomes, solutions and Markov kernels of an SCM can be obtained in this way. For example, mixtures of solutions are also solutions, but not all mixtures can be obtained as the push-forward through a solution function.

**Example 6.2.16.** *For an SCM with endogenous real variables  $X_1, X_2$  and structural equations*

$$\begin{cases} X_1 &= X_2, \\ X_2 &= X_1^3, \end{cases}$$

*any real-valued random variable  $Y$  for which  $P(Y \in \{-1, 0, 1\}) = 1$  provides a solution  $(X_1, X_2) := (Y, Y)$ . This includes all mixtures over the three possible states  $(-1, -1)$ ,  $(0, 0)$ ,  $(1, 1)$ , which form a two-dimensional convex space. However, it has only three solution functions (mapping  $*$  to either  $(-1, -1)$ ,  $(0, 0)$  or  $(1, 1)$ ). Therefore, only three solutions can be constructed from a solution function as in Remark 6.2.15 (namely the extreme points of the convex space).*

One can also use the solution function to sample from the corresponding Markov kernel of an SCM.

**Remark 6.2.17.** *Given a solution function  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  of an SCM  $M = (J, V, W, \mathcal{X}, P, f)$ , we can sample from its corresponding Markov kernel in the following way:*

1. *Input  $x_J$ ;*
2. *For  $w \in W$ : sample  $x_w \sim P(X_w)$ ;*
3. *Calculate  $x_V := g(x_J, x_W)$ ;*
4. *Output  $(x_J, x_V, x_W)$ .*

We can think of this as modeling some data-generating process. If the SCM admits multiple solution functions, this sampler depends explicitly on the choice of the solution function. So one way to think about SCMs that admit multiple solution functions is that they are *incomplete* models of a data-generating process.

Let us now, as a more extensive example, formalize the chocolate-Nobel prize example discussed in Section 1.1 as different SCMs according to some of the causal hypotheses.

**Example 6.2.18.** *For a given country, consider two real-valued variables: annual chocolate consumption in kilograms per capita ( $C$ ), and the number of Nobel prize winners per year per capita ( $N$ ). We can consider the following linear SCM families  $M_\theta = (J, V, W, \mathcal{X}, P_\theta, f_\theta)$ .*

1.  *$N$  causes  $C$ : (“Nobel prizes are celebrated with massive chocolate feasts”)*  
 $J = \{N\}$ ,  $V = \{C\}$ ,  $W = \emptyset$ ,  $\theta = (\alpha, \beta) \in \mathbb{R}^2$ ,  $\mathcal{X}_N = \mathbb{R}$ ,  $\mathcal{X}_C = \mathbb{R}$ ,  $f_{\alpha, \beta} : (x_N, x_C) \mapsto \alpha + \beta x_N$ . *It has structural equations*

$$X_C^{\text{do}(x_N)} = \alpha + \beta x_N.$$

*For a given parameter  $\theta$ , the SCM  $M_\theta$  has a unique solution function,  $g_{\alpha, \beta} : x_N \mapsto \alpha + \beta x_N$ , and unique potential outcomes of the form  $X_C^{\text{do}(x_N)} = \alpha + \beta x_N$ .*

2.  $C$  causes  $N$ : (“chocolate contains brain enhancing chemicals”)  
 $J = \{C\}$ ,  $V = \{N\}$ ,  $W = \emptyset$ ,  $\theta = (\gamma, \delta) \in \mathbb{R}^2$ ,  $\mathcal{X}_C = \mathbb{R}$ ,  $\mathcal{X}_N = \mathbb{R}$ ,  $f_{\gamma, \delta} : (x_C, x_N) \mapsto \gamma + \delta x_C$ . It has structural equations

$$X_N^{\text{do}(x_C)} = \gamma + \delta x_C.$$

For a given parameter  $\theta$ , the SCM  $M_\theta$  has a unique solution function,  $g_{\gamma, \delta} : x_C \mapsto \gamma + \delta x_C$ , and unique potential outcomes of the form  $X_N^{\text{do}(x_C)} = \gamma + \delta x_C$ .

3.  $W$  causes  $C$  and  $N$ , version 1: (“inhabitants of wealthy countries eat more chocolate and conduct more scientific research”)

$J = \{W\}$ ,  $V = \{C, N\}$ ,  $W = \emptyset$ ,  $\theta = (\alpha, \beta, \gamma, \delta) \in \mathbb{R}^4$ ,  $\mathcal{X}_W = \mathbb{R}$ ,  $\mathcal{X}_C = \mathbb{R}$ ,  $\mathcal{X}_N = \mathbb{R}$ ,  $f_\theta : (x_W, x_C, x_N) \mapsto (\alpha + \beta x_W, \gamma + \delta x_W) \in \mathcal{X}_C \times \mathcal{X}_N$ . It has structural equations

$$\begin{aligned} X_C^{\text{do}(x_W)} &= \alpha + \beta x_W, \\ X_N^{\text{do}(x_W)} &= \gamma + \delta x_W. \end{aligned}$$

For a given parameter  $\theta$ , the SCM  $M_\theta$  has unique solution function  $g_\theta : x_W \mapsto (\alpha + \beta x_W, \gamma + \delta x_W) \in \mathcal{X}_C \times \mathcal{X}_N$  and it has unique potential outcomes of the form  $X_{C,N}^{\text{do}(x_W)} = (\alpha + \beta x_W, \gamma + \delta x_W)$ .

4.  $W$  causes  $C$  and  $N$ , version 2: (“similar to version 1, but now the probability distribution of wealth is modeled”):

$J = \emptyset$ ,  $V = \{C, N\}$ ,  $W = \{W\}$ ,  $\theta = (\alpha, \beta, \gamma, \delta, \sigma) \in \mathbb{R}^4 \times [0, \infty)$ ,  $\mathcal{X}_W = \mathbb{R}$ ,  $\mathcal{X}_C = \mathbb{R}$ ,  $\mathcal{X}_N = \mathbb{R}$ ,  $f_\theta : (x_C, x_N, x_W) \mapsto (\alpha + \beta x_W, \gamma + \delta x_W) \in \mathcal{X}_C \times \mathcal{X}_N$ ,  $P_\theta = \mathcal{N}(0, \sigma^2)$ . It has structural equations

$$\begin{aligned} X_C &= \alpha + \beta X_W, \\ X_N &= \gamma + \delta X_W. \end{aligned}$$

For a given parameter  $\theta$ , the SCM  $M_\theta$  has unique solution function  $g_\theta : x_W \mapsto (\alpha + \beta x_W, \gamma + \delta x_W) \in \mathcal{X}_C \times \mathcal{X}_N$  and it has outcomes of the form  $X_{C,N} = (\alpha + \beta X_W, \gamma + \delta X_W)$  for some random variable  $X_W \sim \mathcal{N}(0, \sigma^2)$ .

### 6.3. Interventions

The reason that equations (38) are called *structural* is that one cannot simply rewrite them in the way one is used to when solving a set of equations without changing the causal semantics of the model. This can be formalized by defining how interventions affect an SCM.

In this section we define interventions as *operations* on SCMs that map a given SCM and an intervention target (and optionally, an intervention value or distribution) to an intervened SCM. The operation may change the variable types. We will consider four intervention types: three variants of a hard intervention, and soft interventions. The three hard intervention variants differ in what type of variables the intervened

variables become: endogenous variables, exogenous random variables, or exogenous input variables. As there are many ways in the real world to intervene on (or “to perturb”, or simply “to change”) a given system, this is only the tip of an iceberg of how one could formalize such interventions.<sup>33</sup>

### 6.3.1. Hard interventions

We start with hard interventions that turn all intervened variables into endogenous variables with specified values, overriding the default causal mechanisms that determined their values before the intervention was performed.

**Definition 6.3.1** (Hard intervention with specified target values). *Given an SCM  $M = (J, V, W, \mathcal{X}, P, f)$ , an intervention target  $T \subseteq J \cup V \cup W$  and an intervention value  $\xi_T \in \mathcal{X}_T$ , we define the intervened SCM*

$$M_{\text{do}(X_T=\xi_T)} := (J \setminus T, V \cup T, W \setminus T, \mathcal{X}, P_{W \setminus T}, (f_{V \setminus T}, \xi_T)).$$

More explicitly, the components of the intervened causal mechanism  $\tilde{f} : \mathcal{X} \rightarrow \mathcal{X}_{V \cup T}$  are given by:

$$\tilde{f}_j(x) = \begin{cases} \xi_j & j \in T \\ f_j(x) & j \in V \setminus T, \end{cases}$$

for  $j \in V \cup T$ , and the intervened exogenous distribution is obtained by marginalizing:

$$P_{W \setminus T} = \bigotimes_{w \in W \setminus T} P_w.$$

This replaces the targeted exogenous variables by endogenous variables and adds structural equations to set their values as specified, replaces the existing structural equations of the form  $X_j^{\text{do}(x_j)} = f_j(x_J, X_V^{\text{do}(x_j)}, X_W^{\text{do}(x_j)})$  for  $j \in T \cap V$  to structural equations of the simple form  $X_j^{\text{do}(x_{J \setminus T})} = \xi_j$ , and leaves the other structural equations invariant. This operation “endogenizes” exogenous input variables and exogenous random variables, reflecting that the intervened model now specifies their values as prescribed by the hard intervention. The values of the other endogenous variables are still determined by their original causal mechanisms.

**Example 6.3.2.** *A hard intervention  $\text{do}(X_N = \xi_N)$  changes the SCM “ $W$  causes  $C$  and  $N$ , version 2” from Example 6.2.18 into the SCM with  $J = \emptyset$ ,  $V = \{C, N\}$ ,  $W = \{W\}$ ,  $\theta = (\alpha, \beta, \gamma, \delta, \sigma) \in \mathbb{R}^4 \times [0, \infty)$ ,  $\mathcal{X}_W = \mathbb{R}$ ,  $\mathcal{X}_C = \mathbb{R}$ ,  $\mathcal{X}_N = \mathbb{R}$ ,  $\tilde{f}_\theta : (x_C, x_N, x_W) \mapsto (\alpha + \beta x_W, \xi_N) \in \mathcal{X}_C \times \mathcal{X}_N$ ,  $P_\theta = \mathcal{N}(0, \sigma^2)$ . It has structural equations*

$$\begin{aligned} X_C &= \alpha + \beta X_W, \\ X_N &= \xi_N. \end{aligned}$$

<sup>33</sup>We do not introduce node-splitting interventions for SCMs, as it is not entirely obvious to us how to do this in a sensible way; a possibility might be to generalize these to scc-splitting interventions.

Another common variant of hard interventions are stochastic hard interventions, where the intervention values are drawn independently from a specified (independent) distribution.

**Definition 6.3.3** (Stochastic hard intervention). *Given an SCM  $M = (J, V, W, \mathcal{X}, P, f)$ , an intervention target  $T \subseteq J \cup V \cup W$  and an intervention target distribution  $Q_T \in \bigotimes_{t \in T} \mathcal{P}(\mathcal{X}_t)$ , we define the intervened SCM*

$$M_{\text{do}(X_T \sim Q_T)} := (J \setminus T, V \setminus T, W \cup T, \mathcal{X}, P_{W \setminus T} \otimes Q_T, f_{V \setminus T}).$$

More explicitly, the intervened exogenous distribution is given by

$$P_{W \setminus T} \otimes Q_T = \left[ \bigotimes_{w \in W \setminus T} P_w \right] \otimes \left[ \bigotimes_{t \in T} Q_t \right].$$

Intuitively, this assigns random values to the intervention target variables by sampling from an independent and factorizing intervention distribution  $Q_T$ , thereby turning the targeted variables into exogenous random variables. A hard intervention on an exogenous input variable turning it into an exogenous random variable can be interpreted as “imposing a distribution” on the exogenous input variable. For example, if treatment is considered an exogenous input variable (the model does not specify how treatment is determined by the physician for each patient), and we then intervene to let treatment be determined by a coin flip instead (when setting up an RCT), we are imposing a distribution on the treatment variable.

**Example 6.3.4.** *The SCM “ $W$  causes  $C$  and  $N$ , version 2” from Example 6.2.18 is obtained from a stochastic intervention on the SCM “ $W$  causes  $C$  and  $N$ , version 1” from that example.*

The third variant of hard interventions only specifies the intervention targets, but makes no assertions about the intervention values (not even their distribution).

**Definition 6.3.5** (Hard intervention with unspecified value). *Given an SCM  $M = (J, V, W, \mathcal{X}, P, f)$  and an intervention target  $T \subseteq J \cup V \cup W$ , we define the intervened SCM*

$$M_{\text{do}(T)} := (J \cup T, V \setminus T, W \setminus T, \mathcal{X}, P_{W \setminus T}, f_{V \setminus T}).$$

Intuitively, this operation replaces endogenous variables and exogenous random variables with exogenous input variables. The intervened model no longer specifies the causal mechanisms that determine the values of these variables, but instead treats them as exogenous inputs that are independent of the (remaining) exogenous random variables in the model. This reflects that after this hard intervention, the values for these variables are no longer determined by the system, but are set externally (e.g., by the experimenter performing the intervention) to values chosen independently of the values of the exogenous random variables, while the values of the other endogenous variables are still determined by their original causal mechanisms.

**Example 6.3.6.** A hard intervention  $\text{do}(N)$  changes the SCM “ $W$  causes  $C$  and  $N$ , version 2” from Example 6.2.18 into the intervened SCM with:  $J = \{N\}$ ,  $V = \{C\}$ ,  $W = \{W\}$ ,  $\theta = (\alpha, \beta, \gamma, \delta, \sigma) \in \mathbb{R}^4 \times [0, \infty)$ ,  $\mathcal{X}_W = \mathbb{R}$ ,  $\mathcal{X}_C = \mathbb{R}$ ,  $\mathcal{X}_N = \mathbb{R}$ ,  $\tilde{f}_\theta : (x_N, x_C, x_W) \mapsto \alpha + \beta x_W$ ,  $P_\theta = \mathcal{N}(0, \sigma^2)$ . It has structural equation

$$X_C^{\text{do}(x_N)} = \alpha + \beta X_W^{\text{do}(x_N)}.$$

For a given parameter  $\theta$ , its solution function is unique and given by  $\tilde{g}_\theta : (x_N, x_W) \mapsto \alpha + \beta x_W$ . It has potential outcomes of the form  $X_C^{\text{do}(x_N)} = \alpha + \beta X_W^{\text{do}(x_N)}$  for some random variable  $X_W^{\text{do}(x_N)} \sim \mathcal{N}(0, \sigma^2)$ .

Summarizing, we have now seen three different ways of representing hard interventions, which are all ‘hard’ in the sense that they completely override the default causal mechanisms of their endogenous targets, so that their values are no longer determined by those of other endogenous variables. The three variants differ in how we decide to model the intervened variables: as exogenous inputs, as exogenous random variables, or as endogenous variables with a constant value.<sup>34</sup>

**Proposition 6.3.7.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. Hard interventions  $\text{do}(T_1 \dots)$ ,  $\text{do}(T_2 \dots)$  with disjoint targets  $T_1, T_2 \subseteq J \cup W \cup V$  (of any of the three variants) commute:

$$(M_{\text{do}(T_1 \dots)})_{\text{do}(T_2 \dots)} = (M_{\text{do}(T_2 \dots)})_{\text{do}(T_1 \dots)}.$$

*Proof.* This follows by writing out the definitions and checking commutativity of the operations performed on the various components of the SCM tuple one-by-one.  $\square$

### 6.3.2. Soft interventions (mechanism changes)

Soft interventions replace the causal mechanism of an endogenous variable by another causal mechanism (also known as mechanism changes), or the exogenous distribution of an exogenous random variable by another distribution.

**Definition 6.3.8** (Soft intervention). Given an SCM  $M = (J, V, W, \mathcal{X}, P, f)$  and an intervention target  $T \subseteq V \cup W$ , a soft intervention on  $M$  targeting  $T$  yields an intervened SCM of the form

$$\tilde{M} := (J, V, W, \mathcal{X}, \tilde{P}, \tilde{f})$$

where

$$\tilde{P} = \left( \bigotimes_{w \in W \setminus T} P_w \right) \otimes \left( \bigotimes_{w \in T \cap W} \tilde{P}_w \right)$$

and

$$\tilde{f}(x) = \begin{cases} f_v(x) & v \in V \setminus T \\ \tilde{f}_v(x) & v \in T \cap V \end{cases}$$

for some distributions  $\tilde{P}_w \in \mathcal{P}(\mathcal{X}_w)$  and measurable functions  $\tilde{f}_v : \mathcal{X} \rightarrow \mathcal{X}_v$ .

<sup>34</sup>Since most accounts on SCMs do not provide for the possibility of exogenous input variables, the two variants most often seen in the literature are the latter two.

### 6.3.3. Intervention variables

It is often convenient to combine an SCM and one or more intervened versions of the SCM together in a single SCM. This can be done by introducing *intervention variables* that indicate whether, and possibly encode how, an intervention is performed.

We give no rigorous definition—deciding on how to make use of intervention variables is a modeling issue—but provide a brief discussion.

Often, intervention variables can be considered as exogenous. For instance, if their values (which intervention is performed, and how) are determined earlier in time than the values of all endogenous variables, then the endogenous variables cannot cause the intervention variables. An example is a scientific experiment in which the experimenter prepares the system in some experimental condition before performing measurements. This experimental condition then cannot be influenced by the measured state of the system. A concrete instance of such an intervention variable is the coin flip that determines whether a subject enters the treatment or control group in a randomized controlled trial.

There is no need for an intervention variable to be binary or discrete. An example of a continuous intervention variable is the dose of the drug used for treatment. The dose of a drug used for treatment often quantitatively affects the outcome. However, if there is no explicit and careful randomization, one cannot always assume that there is no common cause of dose and response. For example, if physicians decide to not treat patients with terminal cancer because of the strong side effects of the treatment, then the stage of the cancer is a common cause of dose and response.

Therefore, for intervention variables (like for any other variable that is modeled) one needs to take care to decide whether they are modeled as exogenous input variables, exogenous random variables, or endogenous variables.

## 6.4. Composition and decomposition

If we think about an SCM as modeling a “system”, then we also obtain a model for any “subsystem” in the following way.

**Definition 6.4.1** (Taking a submodel of an SCM). *Given an SCM  $M = (J, V, W, \mathcal{X}, P, f)$  and a subset  $V' \subseteq V$  of its endogenous variables, we define its submodel on  $V'$  as the SCM*

$$M_{[V']} := (J \cup (V \setminus V'), V', W, \mathcal{X}, P, f_{V'}).$$

*Note that this is just the intervened SCM  $M_{\text{do}(V \setminus V')}$ .*

Given two SCMs, such that some of the variables of one SCM can be used as (part of the) exogenous input of the other, and possibly vice versa, we can compose them into a single SCM.<sup>35</sup>

**Definition 6.4.2** (Composing two SCMs). *Given two SCMs  $M_1 = (J_1, V_1, W_1, \mathcal{X}_1, P_1, f_1)$  and  $M_2 = (J_2, V_2, W_2, \mathcal{X}_2, P_2, f_2)$  and two subsets  $C_1 \subseteq J_1 \cap (V_2 \cup W_2)$  and  $C_2 \subseteq (V_1 \cup W_1) \cap J_2$  satisfying the following conditions:*

<sup>35</sup>In practice, one may have to relabel the variables first before one can perform this composition operation, by changing the index sets of the SCMs to ensure that the right variables become coupled.



1. for all  $c \in C_1$ ,  $(\mathcal{X}_1)_c = (\mathcal{X}_2)_c$ ,
2. for all  $c \in C_2$ ,  $(\mathcal{X}_1)_c = (\mathcal{X}_2)_c$ ,
3.  $(J_1 \setminus C_1) \cap (J_2 \setminus C_2) = \emptyset$ ,
4.  $V_1 \cap V_2 = \emptyset$ ,
5.  $W_1 \cap W_2 = \emptyset$ ,

we define the composed SCM

$$M_{12} := ((J_1 \setminus C_1) \dot{\cup} (J_2 \setminus C_2), V_1 \dot{\cup} V_2, W_1 \dot{\cup} W_2, \mathcal{X}_{12}, P_{12}),$$

where

$$\begin{aligned} \mathcal{X}_{12} &:= (\mathcal{X}_1)_{(J_1 \setminus C_1) \dot{\cup} V_1 \dot{\cup} W_1} \times (\mathcal{X}_2)_{(J_2 \setminus C_2) \dot{\cup} V_2 \dot{\cup} W_2}, \\ P_{12} &:= P_1 \otimes P_2, \\ f_{12} &:= (f_1, f_2). \end{aligned}$$

The special case  $C_1 = \emptyset$  will also be denoted as  $M_{2\circ 1}$ , and likewise the special case  $C_2 = \emptyset$  will be denoted as  $M_{1\circ 2}$ .

One can consider the decomposition (taking submodels) as ‘dual to’ the composition (combining submodels). These operations formalize the notion of *modularity*, that is, an SCM can be thought of modeling a system consisting of interacting components by modeling for each component separately how it interacts with the other components. The submodels on individual variables correspond with ‘atomic’ subsystems that cannot be (or won’t be) decomposed into smaller parts.

## 6.5. Unique solvability and simple SCMs

**Definition 6.5.1.** An SCM  $M = (J, V, W, \mathcal{X}, P, f)$  is called uniquely solvable if it has a unique solution function. This means two things: (i) it has a solution function  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$ , i.e., a measurable function that satisfies

$$\forall x_W \in \mathcal{X}_W \forall x_J \in \mathcal{X}_J : g(x_J, x_W) = f(x_J, g(x_J, x_W), x_W); \quad (41)$$

(ii) all its solution functions must equal  $g$ , i.e., for any solution function  $\tilde{g} : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$ , we have that for all  $x_W \in \mathcal{X}_W$ , for all  $x_J \in \mathcal{X}_J$ :

$$g(x_J, x_W) = \tilde{g}(x_J, x_W).$$

The following result gives two useful properties of unique solvability. The first provides an equivalent formulation that makes it easier to check whether an SCM is uniquely solvable, the second provides an important consequence regarding the Markov kernels of the SCM.

**Theorem 6.5.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM.*

1.  *$M$  is uniquely solvable if and only if for all  $x_W \in \mathcal{X}_W$ , for all  $x_J \in \mathcal{X}_J$ , the equation*

$$x_V = f(x_J, x_V, x_W)$$

*has a unique solution for  $x_V \in \mathcal{X}_V$ .*

2. *If  $M$  is uniquely solvable, it has a unique Markov kernel*

$$P_M(X_V, X_W \mid \text{do}(X_J)) = (g, \text{id}_{\mathcal{X}_W})_* \left( \left( \bigotimes_{w \in W} P_w(X_w) \right) \otimes \delta(X_J \mid X_J) \right)$$

*where  $g$  is the unique solution function of  $M$ . The distribution of any potential outcome  $X_{V \cup W}^{\text{do}(x_J)}$  for input  $x_J \in X_J$  is given by  $P_M(X_V, X_W \mid \text{do}(X_J = x_J))$ .*

*Proof.* 1. “ $\implies$ ”: Let  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  be a solution function for  $M$ . Let  $\xi_W \in \mathcal{X}_W$ ,  $\xi_J \in \mathcal{X}_J$ . The equation  $x_V = f(x_V, \xi_J, \xi_W)$  does have a solution for  $x_V$  (indeed,  $g(\xi_J, \xi_W)$  is such a solution). Suppose its solution is not unique, i.e., it has another solution  $\tilde{x}_v \neq g(\xi_J, \xi_W)$ . Consider the modified function

$$\tilde{g} : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V : (x_J, x_W) \mapsto \begin{cases} \tilde{x}_v & x_J = \xi_J \wedge x_W = \xi_W \\ g(x_J, x_W) & \text{otherwise} \end{cases}$$

$\tilde{g}$  is measurable (because  $g$  is measurable), satisfies (41), and hence it provides a solution function. However,  $\tilde{g} \neq g$ , which contradicts the assumed unique solvability.

“ $\impliedby$ ”: This boils down to proving that the measurability of  $f$  and the uniqueness of the solutions implies the measurability of the solution function. We exploit that we are dealing with standard measurable spaces. Define the function  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  by letting  $g(x_J, x_W)$  be the (unique) solution  $x_V$  of the equation  $x_V = f(x_V, x_J, x_W)$ , with  $x_V \in \mathcal{X}_V$ . The graph of this function is

$$\text{graph}(g) = \{(x_J, x_W, x_V) \in \mathcal{X}_J \times \mathcal{X}_W \times \mathcal{X}_V : x_V = g(x_J, x_W)\}.$$

By assumption, we have that  $x_V = g(x_J, x_W) \iff x_V = f(x_V, x_J, x_W)$  for all  $x \in \mathcal{X}$ . Hence,

$$\text{graph}(g) = \{(x_J, x_W, x_V) \in \mathcal{X}_J \times \mathcal{X}_W \times \mathcal{X}_V : x_V = f(x_V, x_J, x_W)\}.$$

Defining the function

$$h : \mathcal{X}_J \times \mathcal{X}_W \times \mathcal{X}_V \rightarrow \mathcal{X}_V \times \mathcal{X}_V : (x_J, x_W, x_V) \mapsto (x_V, f(x_V, x_J, x_W))$$

and the diagonal  $\Delta = \{(x_V, x_V) : x_V \in \mathcal{X}_V\} \subseteq \mathcal{X}_V^2$ , this shows that

$$\text{graph}(g) = h^{-1}(\Delta).$$

Since  $h$  is measurable and  $\Delta$  is a measurable set (because  $\mathcal{X}_V$  is Hausdorff),  $h^{-1}(\Delta)$  is a measurable set. By [Kec95, 14.12], because all spaces are (isomorphic to) Borel spaces, the fact that  $\text{graph}(g)$  is a measurable set implies that  $g$  is a measurable function. Hence  $g$  is a solution function. The unique solvability of  $M$  follows since this is the only possible solution function.

2. Assume  $M$  to be uniquely solvable and let  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  be its unique solution function. The push-forward  $(g, \text{id}_{\mathcal{X}_W})_*(P)$  of the exogenous distribution  $P$  of  $M$  (interpreted as a constant Markov kernel from  $\mathcal{X}_J$  to  $\mathcal{X}_W$ ) provides a Markov kernel for  $M$ .

Let  $K(X_V, X_W | X_J)$  denote any Markov kernel for  $M$  corresponding to a solution  $X : \mathcal{U} \times \mathcal{X}_J \rightarrow \mathcal{X}$ . Pick any  $x_J \in \mathcal{X}_J$ . We have to show that  $K(X_V, X_W | X_J = x_J)$  is a unique distribution. The distribution  $K(X_V, X_W | X_J = x_J)$  is that of  $X_{V \cup W}^{\text{do}(x_J)}$ , which is a potential outcome of  $M$ .

Since  $X_{V \cup W}^{\text{do}(x_J)}$  is a potential outcome of  $M$ , we have that  $X_W^{\text{do}(x_J)} \sim P$  and

$$X_V^{\text{do}(x_J)} = f(x_J, X_V^{\text{do}(x_J)}, X_W^{\text{do}(x_J)}) \quad \text{a.s.}$$

This implies that

$$X_V^{\text{do}(x_J)} = g(x_J, X_W^{\text{do}(x_J)}) \quad \text{a.s.}$$

By modifying the random variable  $X_{V \cup W}^{\text{do}(x_J)}$  on a null set, we can obtain a random variable  $Y_{V \cup W}^{\text{do}(x_J)}$  such that we get equality everywhere:

$$Y_V^{\text{do}(x_J)} = g(x_J, Y_W^{\text{do}(x_J)}).$$

This shows that the distribution of  $Y_{V \cup W}^{\text{do}(x_J)}$ , and hence that of  $X_{V \cup W}^{\text{do}(x_J)}$ , is that of the push-forward of the exogenous distribution  $P$  through the unique function  $g_{x_J} : \mathcal{X}_W \rightarrow \mathcal{X}_V : x_W \mapsto (g(x_J, x_W), x_W)$ . The latter push-forward distribution is unique. □

The first equivalence states that one gets the *measurability* of the solution function for free from the measurability of the causal mechanism  $f$  and the uniqueness of the solutions of the structural equations. Note that for the special case of no exogenous input variables ( $J = \emptyset$ ), the above shows that uniquely solvable SCMs induce a unique distribution.

For SCMs whose causal mechanism is linear in terms of the endogenous variables, we can give a sufficient condition for unique solvability, and explicitly write down the form of their solution function.

**Definition 6.5.3.** *An SCM  $M = (J, V, W, \mathcal{X}, P, f)$  is called linear if all endogenous variables are real-valued (i.e.,  $\mathcal{X}_v = \mathbb{R}$  for  $v \in V$ ), and each component of the causal*

mechanism is an affine combination of endogenous variables with coefficients that may depend on exogenous variables, i.e., of the form

$$f_v(x) = \sum_{u \in V} B_{vu}(x_J, x_W) x_u + c_v(x_J, x_W),$$

where  $B(x_J, x_W) \in \mathbb{R}^{V \times V}$  is a family of matrices (or: a matrix-valued function  $\mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathbb{R}^{V \times V}$ ) and  $c(x_J, x_W) \in \mathbb{R}^V$  is a family of real-valued offsets (or: a vector-valued function  $\mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathbb{R}^V$ ).

For linear SCMs, unique solvability is equivalent to the invertability of a certain matrix.

**Proposition 6.5.4.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a linear SCM. Then, for any set  $L \subseteq V$ , the submodel  $M_{[L]}$  is uniquely solvable if and only if the matrices  $I_L - B_{LL}(x_J, x_W)$  are invertible for all  $x_J \in \mathcal{X}_J, x_W \in \mathcal{X}_W$ , where  $I_L \in \mathbb{R}^{L \times L}$  the identity matrix and  $B_{LL}(x_J, x_W) \in \mathbb{R}^{L \times L}$  the submatrix of  $B(x_J, x_W)$ . Its unique solution function is then:*

$$\begin{aligned} g : \mathbb{R}^{V \setminus L} \times \mathbb{R}^{J \cup W} &\rightarrow \mathbb{R}^L \\ : (x_{V \setminus L}, x_J, x_W) &\mapsto (I_L - B_{LL}(x_J, x_W))^{-1} (B_{L, V \setminus L}(x_J, x_W) x_{V \setminus L} + c_L(x_J, x_W)). \end{aligned}$$

If an SCM is uniquely solvable, this does not necessarily mean that it is still uniquely solvable after performing some intervention. To avoid the complications introduced in case solutions are absent, or present but not unique, we will henceforth make strong assumptions regarding the existence and uniqueness of solutions. We (mostly) restrict our attention to a subclass of SCMs that we refer to as *simple SCMs*, which are SCMs that are uniquely solvable and remain so after any hard intervention:

**Definition 6.5.5.** *An SCM  $M = (J, V, W, \mathcal{X}, P, f)$  is called simple if the intervened SCM  $M_{\text{do}(T)}$  is uniquely solvable for all  $T \subseteq J \cup V \cup W$ .*

Note that this includes unique solvability of  $M$  itself for  $T = \emptyset$ .

**Remark 6.5.6.** *It suffices to check whether  $M_{\text{do}(T)}$  is uniquely solvable for all  $T \subseteq V$ , since interventions on exogenous variables turn them into exogenous input variables, leaving the solution function invariant.*

**Example 6.5.7.** *Consider an SCM with structural equations*

$$\begin{aligned} X_1 &= W_1 \\ X_2 &= W_2 \\ X_3 &= X_1 X_4 + W_3 \\ X_4 &= X_2 X_3 + W_4 \end{aligned}$$

where the  $X$ 's are considered real-valued endogenous variables and the  $W$ 's exogenous variables with domains  $(-1, 1) \subset \mathbb{R}$ . We can solve the system of structural equations for

$X$  in terms of  $W$ :

$$\begin{aligned} X_1 &= W_1 \\ X_2 &= W_2 \\ X_3 &= \frac{W_3 + W_1W_4}{1 - W_1W_2} \\ X_4 &= \frac{W_2W_3 + W_4}{1 - W_2W_1} \end{aligned}$$

Similarly, we can take any subset of the structural equations and solve it for the variables appearing on the l.h.s. of the equations in the subset, and obtain a unique solution. For example, only solving the structural equations for  $X_3$  and  $X_4$ , we obtain:

$$\begin{aligned} X_3 &= \frac{W_3 + W_1W_4}{1 - W_1X_2} \\ X_4 &= \frac{W_2W_3 + W_4}{1 - W_2X_1} \end{aligned}$$

where the variables  $X_1$  and  $X_2$  are now considered as exogenous input variables (instead of endogenous variables). Hence, any subset of the structural equations has a unique solution for the variables appearing on the l.h.s. in terms of the remaining ones on the r.h.s., which means that this SCM is simple.

**Corollary 6.5.8.** *If  $M$  is simple, then its Markov kernel  $P_M(X_V, X_W \mid \text{do}(X_J))$  and all intervened Markov kernels*

$$P_M(X_{(V \cup W) \setminus T} \mid \text{do}(X_J), \text{do}(X_T)) := P_{M_{\text{do}(T)}}(X_{(V \cup W) \setminus T} \mid \text{do}(X_{J \cup T}))$$

for  $T \subseteq V \cup W$  exist and are unique.

*Proof.* This follows from Theorem 6.5.2. □

The two notations for intervened Markov kernels of simple SCMs will be used interchangeably.

The fact that these SCMs are relatively simple to deal with (because we do not have to worry about non-existence or non-uniqueness of solutions) motivated their name. Even though simplicity is a strong assumption, the class of simple SCMs is more expressive than the class of L-CBNs since it allows to model causal cycles, to some extent.

## 6.6. Equivalences

Equivalence relations are ubiquitous in mathematics. They capture the notion that mathematical objects can be “equivalent” from some point of view.

**Definition 6.6.1.** *Let  $Z$  be a set and  $R \subseteq Z^2$  be a relation on  $Z$  (i.e., a subset of ordered pairs of  $Z$ ).  $R$  is called an equivalence relation if*

1.  $R$  is reflexive:  $(a, a) \in R$  for all  $a \in Z$ ;
2.  $R$  is symmetric:  $(a, b) \in R \iff (b, a) \in R$  for all  $a, b \in Z$ ;
3.  $R$  is transitive: if  $(a, b) \in R$  and  $(b, c) \in R$  then  $(a, c) \in R$  for all  $a, b, c \in Z$ .

In this section, we will discuss several important equivalence relations between SCMs.

**Definition 6.6.2.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  and  $\tilde{M} = (J, V, W, \mathcal{X}, P, \tilde{f})$  be two SCMs that may differ only in terms of their causal mechanism. We say that  $M$  is equivalent to  $\tilde{M}$  and write  $M \equiv \tilde{M}$  if for each  $v \in V$ ,

$$\forall x \in \mathcal{X} : \quad x_v = f_v(x_J, x_V, x_W) \iff x_v = \tilde{f}_v(x_J, x_V, x_W).$$

In words,  $M$  and  $\tilde{M}$  are considered equivalent if they at most differ in terms of their causal mechanisms, yet each of their structural equations has the same solutions.

Equivalent SCMs are indeed equivalent for many purposes:

**Proposition 6.6.3.** 1. *Equivalence is preserved by hard interventions:*

if  $M \equiv \tilde{M}$  and  $T \subseteq V \cup W \cup J$ , then  $M_{\text{do}(T\dots)} \equiv \tilde{M}_{\text{do}(T\dots)}$ ;

2. *Unique solvability is an invariant of equivalence:*

if  $M \equiv \tilde{M}$  then  $M$  is uniquely solvable if and only if  $\tilde{M}$  is uniquely solvable;

3. *Simplicity is an invariant of equivalence:*

if  $M \equiv \tilde{M}$  then  $M$  is simple if and only if  $\tilde{M}$  is simple;

4. *Equivalent SCMs have the same solution functions, solutions, (potential) outcomes, and Markov kernels.*

*Proof.* Exercise for the reader. □

**Example 6.6.4.** *The SCM with the single structural equation*

$$X = -X^3 + X + W$$

where  $X$  is endogenous, and  $W$  is exogenous, is equivalent to the SCM obtained by replacing the structural equation with

$$X = \sqrt[3]{W}.$$

Note that  $X$  no longer appears on the r.h.s..

This notion of equivalence is rather strong.

**Example 6.6.5.** *The SCM with endogenous variable  $X \in \mathbb{R}$  and exogenous random variable  $W \in \mathbb{R}$  and structural equation*

$$X = WX + c$$

*is not equivalent to the SCM obtained by replacing the structural equation with*

$$X = \begin{cases} \xi & W = 1 \\ \frac{c}{1-W} & W \neq 1 \end{cases}$$

*no matter how we choose  $\xi$ , even when  $P(W = 1) = 0$ . If one does not model interventions on exogenous random variables, weaker notions of equivalence can be used that replace the “for all” quantifier over values of exogenous random variables by the “for almost all” quantifier (see [BFPM21]).*

We often make use of other notions of equivalence as well. For simplicity of exposition, we provide the definitions only for simple SCMs (the general definitions are provided in [BFPM21] for SCMs without exogenous input variables).

**Definition 6.6.6.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  and  $\tilde{M} = (\tilde{J}, \tilde{V}, \tilde{W}, \tilde{\mathcal{X}}, \tilde{P}, \tilde{f})$  be two simple SCMs and  $O \subseteq (V \cup W) \cap (\tilde{V} \cup \tilde{W})$  a subset. We say that:*

1.  *$M$  and  $\tilde{M}$  are observably equivalent w.r.t.  $O$  if  $\mathcal{X}_O = \tilde{\mathcal{X}}_O$ ,  $\mathcal{X}_{J \cap \tilde{J}} = \tilde{\mathcal{X}}_{J \cap \tilde{J}}$  and their marginal Markov kernels coincide:*

$$P_M(X_O \mid \text{do}(X_J)) = P_{\tilde{M}}(X_O \mid \text{do}(X_{\tilde{J}})).$$

*This has to be interpreted as both Markov kernels being a version of a Markov kernel  $\mathcal{X}_{J \cap \tilde{J}} \dashrightarrow \mathcal{X}_O$ , i.e.,  $P_M(X_O \mid \text{do}(X_J))$  must be essentially constant in  $x_{J \setminus \tilde{J}}$  and  $P_{\tilde{M}}(X_O \mid \text{do}(X_{\tilde{J}}))$  must be essentially constant in  $x_{\tilde{J} \setminus J}$ .<sup>36</sup>*

2.  *$M$  and  $\tilde{M}$  are interventionally equivalent w.r.t.  $O$  if for every subset  $T \subseteq O$  the intervened SCMs  $M_{\text{do}(T)}$  and  $\tilde{M}_{\text{do}(T)}$  are observably equivalent w.r.t.  $O \setminus T$ ;*

More generally, one could define interventional equivalence not only with respect to an observed set of variables, but also with respect to a given set of interventions.

One can show the following properties of these equivalences:

**Proposition 6.6.7.** *For simple SCMs  $M, \tilde{M}$  and a subset  $O \subseteq (V \cap \tilde{V}) \cup (W \cap \tilde{W})$ :*

1. *If  $M \equiv \tilde{M}$  then  $M$  and  $\tilde{M}$  are interventionally equivalent w.r.t.  $O$ .*

---

<sup>36</sup>In other words,

$$P_M(X_O \mid \text{do}(X_J)) = P_M(X_O \mid \text{do}(X_{J \cap \tilde{J}})) = P_{\tilde{M}}(X_O \mid \text{do}(X_{J \cap \tilde{J}})) = P_{\tilde{M}}(X_O \mid \text{do}(X_{\tilde{J}})).$$

2. If  $M$  and  $\tilde{M}$  are interventionally equivalent w.r.t.  $O$  then  $M$  and  $\tilde{M}$  are observably equivalent w.r.t.  $O$ .

*Proof.* 1. Suppose  $M \equiv \tilde{M}$ . Then for every  $T \subseteq O$ ,  $M_{\text{do}(T)} \equiv \tilde{M}_{\text{do}(T)}$ . Since equivalent SCMs have the same Markov kernels,  $M_{\text{do}(T)}$  and  $\tilde{M}_{\text{do}(T)}$  are observably equivalent w.r.t.  $O \setminus T$  for every  $T \subseteq O$ .

2. This is trivial (consider the intervention targeting  $T = \emptyset$ ).

□

However, the reverse implications do not hold in general. This expresses that causal modeling is more refined than probabilistic modeling.

**Example 6.6.8** (Observable equivalence does not imply interventional equivalence). Consider the SCM  $M$  with

$$\begin{aligned} N &\sim \mathcal{N}(\mu, \sigma^2) \\ C &= \alpha + \beta N \end{aligned}$$

and the SCM  $\tilde{M}$  with

$$\begin{aligned} C &\sim \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2) \\ N &= \tilde{\alpha} + \tilde{\beta} C \end{aligned}$$

These SCMs are simple and their distributions are respectively

$$P_M(N, C) = \mathcal{N} \left( \begin{pmatrix} \mu \\ \alpha + \beta\mu \end{pmatrix}, \begin{pmatrix} \sigma^2 & \beta\sigma^2 \\ \beta\sigma^2 & \beta^2\sigma^2 \end{pmatrix} \right)$$

and

$$P_{\tilde{M}}(N, C) = \mathcal{N} \left( \begin{pmatrix} \tilde{\alpha} + \tilde{\beta}\tilde{\mu} \\ \tilde{\mu} \end{pmatrix}, \begin{pmatrix} \tilde{\beta}^2\tilde{\sigma}^2 & \tilde{\beta}\tilde{\sigma}^2 \\ \tilde{\beta}\tilde{\sigma}^2 & \tilde{\sigma}^2 \end{pmatrix} \right)$$

For certain parameter choices, they are observably equivalent. However, they are not interventionally equivalent except for very special parameter choices.

## 6.7. Marginalizations

When modeling a system, we sometimes want to “hide” details of a subsystem. The following operation on SCMs that we call “marginalization” is a causal analogue of the marginalization of probability distributions. The computer program analogy of the marginalization operation is to hide details within a subroutine. Intuitively, a marginalization of an SCM over a subset of endogenous variables  $L$  is obtained by first solving a subsystem (the structural equations corresponding to the endogenous variables in  $L$ ) followed by substituting the solution function of the subsystem into the remaining structural equations (corresponding to the endogenous variables in  $V \setminus L$ ).



**Definition 6.7.1.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM and  $L \subseteq V$  such that the submodel  $M_{[L]}$  is uniquely solvable. Let  $\tilde{g}_L : \mathcal{X}_J \times \mathcal{X}_{V \setminus L} \times \mathcal{X}_W \rightarrow \mathcal{X}_L$  be the unique solution function for  $M_{[L]}$ . Then we call  $M_{\setminus L} = (J, V \setminus L, W, \mathcal{X}_J \times \mathcal{X}_{V \setminus L} \times \mathcal{X}_W, P, \tilde{f})$  with

$$\tilde{f}(x_J, x_{V \setminus L}, x_W) = f_{V \setminus L}(x_J, x_{V \setminus L}, \tilde{g}_L(x_J, x_{V \setminus L}, x_W), x_W)$$

the marginalization of  $M$  over  $L$ .

For simple SCMs, marginalizations are obviously defined over any subset  $L \subseteq V$ .

Marginalization preserves unique solvability, as the following lemma shows.

**Lemma 6.7.2.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM and  $L \subseteq V$  such that  $M_{[L]}$  is uniquely solvable. If  $M$  is uniquely solvable then its marginalization  $M_{\setminus L}$  is uniquely solvable, and the unique solution function for  $M_{\setminus L}$  is  $g_{V \setminus L} = \text{pr}_{V \setminus L} \circ g$ , where  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  is the unique solution function for  $M$  and  $\text{pr}_K : \mathcal{X}_V \rightarrow \mathcal{X}_K : x \mapsto x_K$  is the canonical projection on  $K \subseteq V$ .

*Proof.* Let  $K := V \setminus L$  and let  $\tilde{g}_L : \mathcal{X}_{J \cup K} \times \mathcal{X}_W \rightarrow \mathcal{X}_L$  be the unique solution function for  $M_{[L]} = M_{\text{do}(K)}$ . From the properties of the solution functions, we derive that for all  $x \in \mathcal{X}$ ,

$$\begin{aligned} \begin{cases} x_L = g_L(x_J, x_W) \\ x_K = g_K(x_J, x_W) \end{cases} &\iff x_V = g(x_J, x_W) \iff x_V = f(x) \iff \begin{cases} x_L = f_L(x) \\ x_K = f_K(x) \end{cases} \\ \iff \begin{cases} x_L = \tilde{g}_L(x_J, x_K, x_W) \\ x_K = f_K(x_J, x_K, x_L, x_W) \end{cases} &\iff \begin{cases} x_L = \tilde{g}_L(x_J, x_K, x_W) \\ x_K = f_K(x_J, x_K, \tilde{g}_L(x_J, x_K, x_W), x_W) \end{cases} \\ \iff \begin{cases} x_L = \tilde{g}_L(x_J, x_K, x_W) \\ x_K = \tilde{f}(x_J, x_K, x_W) \end{cases} & \end{aligned}$$

where  $\tilde{f}$  is the causal mechanism of  $M_{\setminus L}$ . Hence for all  $x_J \in \mathcal{X}_J, x_W \in \mathcal{X}_W, x_K \in \mathcal{X}_K$ :

$$x_K = g_K(x_J, x_W) \iff x_K = \tilde{f}(x_J, x_K, x_W).$$

Therefore,  $g_K = \text{pr}_K \circ g$  is the unique solution function for  $M_{\setminus L}$ , and the marginalized SCM  $M_{\setminus L}$  is uniquely solvable.  $\square$

**Remark 6.7.3.** This also directly implies that if  $M$  is uniquely solvable and its marginalization  $M_{\setminus L}$  over  $L \subseteq V$  is defined, the Markov kernel of the marginalization is obtained by marginalizing the original Markov kernel:

$$P_{M_{\setminus L}}(X_{V \setminus L}, X_W \mid \text{do}(X_J)) = P_M(X_{V \setminus L}, X_W \mid \text{do}(X_J)).$$

This explains the name ‘marginalization’.

Under certain conditions, hard interventions and marginalization commute (i.e., it does not matter in which order we apply them).

**Proposition 6.7.4.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. For  $L \subseteq V$  such that  $M_{[L]}$  is uniquely solvable, and a hard intervention  $\text{do}(T \dots)$  with target  $T \subseteq J \cup W \cup V$  (of any of the three variants) such that  $L \cap T = \emptyset$ :*

$$(M_{\text{do}(T \dots)})_{\setminus L} = (M_{\setminus L})_{\text{do}(T \dots)}.$$

*Proof.* This follows by writing out the definitions and checking commutativity of the operations performed on the various components of the SCM tuple one-by-one.  $\square$

We will show that for simple SCMs, the marginalization operation preserves the causal semantics on the remaining variables. A key step is the following proposition, which gives conditions under which it does not matter whether we marginalize at once or in steps.

**Proposition 6.7.5.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM and  $L_1, L_2 \subseteq V$  such that  $L_1 \cap L_2 = \emptyset$ . If  $M_{[L_1]}$  is uniquely solvable, and  $(M_{\setminus L_1})_{[L_2]}$  is uniquely solvable, then  $M_{[L_1 \cup L_2]}$  is uniquely solvable, and in that case it does not matter if we first marginalize over  $L_1$  and then  $L_2$ , or both at once, i.e.:*

$$(M_{\setminus L_1})_{\setminus L_2} = M_{\setminus (L_1 \cup L_2)}.$$

*Proof.* Write  $K_1 = V \setminus L_1$ . Let  $\tilde{g} : \mathcal{X}_{J \cup K_1} \times \mathcal{X}_W \rightarrow \mathcal{X}_{L_1}$  be the unique solution function of  $M_{[L_1]}$ , i.e., for all  $x \in \mathcal{X}$ :

$$x_{L_1} = \tilde{g}(x_J, x_{K_1}, x_W) \iff x_{L_1} = f_{L_1}(x).$$

Let  $\tilde{f} : \mathcal{X}_J \times \mathcal{X}_{K_1} \times \mathcal{X}_W \rightarrow \mathcal{X}_{K_1}$  with

$$\tilde{f}(x_J, x_{K_1}, x_W) = f_{K_1}(x_J, x_{K_1}, \tilde{g}(x_J, x_{K_1}, x_W), x_W)$$

be the causal mechanism of the marginal SCM  $M_{\setminus L_1}$ .

If  $(M_{\setminus L_1})_{[L_1 \cup L_2]}$  is uniquely solvable, it has a unique solution function  $\tilde{\tilde{g}} : \mathcal{X}_J \times \mathcal{X}_{V \setminus (L_1 \cup L_2)} \times \mathcal{X}_W \rightarrow \mathcal{X}_{L_2}$ , i.e., for all  $x \in \mathcal{X}$ :

$$x_{L_2} = \tilde{\tilde{g}}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W) \iff x_{L_2} = \tilde{f}_{L_2}(x_J, x_{L_2}, x_{K_1 \setminus L_2}, x_W)$$

Define the function  $h : \mathcal{X}_J \times \mathcal{X}_{V \setminus (L_1 \cup L_2)} \times \mathcal{X}_W \rightarrow \mathcal{X}_{L_1 \cup L_2}$  by

$$\begin{aligned} h_{L_1}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W) &= \tilde{g}(x_J, x_{K_1 \setminus L_2}, \tilde{\tilde{g}}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W), x_W) \\ h_{L_2}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W) &= \tilde{\tilde{g}}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W) \end{aligned}$$

Then for all  $x \in \mathcal{X}$ :

$$\begin{aligned}
& \begin{cases} x_{L_1} &= f_{L_1}(x) \\ x_{L_2} &= f_{L_2}(x) \end{cases} \\
& \iff \begin{cases} x_{L_1} &= \tilde{g}(x_J, x_{K_1}, x_W) \\ x_{L_2} &= f_{L_2}(x_J, x_{K_1}, x_{L_1}, x_W) \end{cases} \\
& \iff \begin{cases} x_{L_1} &= \tilde{g}(x_J, x_{K_1}, x_W) \\ x_{L_2} &= f_{L_2}(x_J, x_{K_1}, \tilde{g}(x_J, x_{K_1}, x_W), x_W) \end{cases} \\
& \iff \begin{cases} x_{L_1} &= \tilde{g}(x_J, x_{K_1 \setminus L_2}, x_{L_2}, x_W) \\ x_{L_2} &= \tilde{\tilde{g}}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W) \end{cases} \\
& \iff \begin{cases} x_{L_1} &= \tilde{g}(x_J, x_{K_1 \setminus L_2}, \tilde{\tilde{g}}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W), x_W) \\ x_{L_2} &= \tilde{\tilde{g}}(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W) \end{cases} \\
& \iff x_{L_1 \cup L_2} = h(x_J, x_{V \setminus (L_1 \cup L_2)}, x_W)
\end{aligned}$$

where in the first equivalence we used the unique solvability of  $M_{[L_1]}$ , in the second equivalence we used substitution, in the third equivalence we used the unique solvability of  $(M_{\setminus L_1})_{[L_2]}$ , in the fourth equivalence we used substitution again, and in the fifth equivalence we used the definition of  $h$ . Therefore,  $h$  is the unique solution function for  $M_{[L_1 \cup L_2]}$ , which must therefore be uniquely solvable. By checking the definition, one concludes that  $(M_{\setminus L_1})_{\setminus L_2} = M_{\setminus (L_1 \cup L_2)}$ .  $\square$

**Proposition 6.7.6.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. For any  $L \subseteq V$ , its marginalization  $M_{\setminus L}$  is also simple.*

*Proof.* Let  $T \subseteq J \cup W \cup V \setminus L$ . By Proposition 6.7.4, the marginalization commutes with the hard intervention:

$$(M_{\setminus L})_{\text{do}(T)} = (M_{\text{do}(T)})_{\setminus L}.$$

Because  $M$  is simple, also  $M_{\text{do}(T)}$  is simple (by Proposition 6.3.7), and in particular it is uniquely solvable. From Lemma 6.7.2 it follows that also its marginalization  $(M_{\text{do}(T)})_{\setminus L}$  is uniquely solvable. This means that  $(M_{\setminus L})_{\text{do}(T)}$  is uniquely solvable. Since this holds for any  $T \subseteq J \cup W \cup V \setminus L$ ,  $M_{\setminus L}$  is simple.  $\square$

**Corollary 6.7.7.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. For  $L_1, L_2 \subseteq V$  with  $L_1 \cap L_2 = \emptyset$ ,*

$$(M_{\setminus L_1})_{\setminus L_2} = (M_{\setminus L_2})_{\setminus L_1} = M_{\setminus (L_1 \cup L_2)}.$$

These commutation relations and compatibilities now allow us to give a straightforward proof that the causal semantics are preserved under marginalization. While this holds generally [BFPM21], we will here only prove this for simple SCMs.

**Theorem 6.7.8.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM,  $L \subseteq V$ , and  $M_{\setminus L}$  its marginalization over  $L$ . Then  $M$  and  $M_{\setminus L}$  are observably and interventionally equivalent w.r.t.  $(V \cup W) \setminus L$ .*

*Proof.* Write  $K = V \setminus L$ . We first show that the marginal Markov kernels  $P_M(X_K, X_W \mid \text{do}(X_J))$  and  $P_{M_{\setminus L}}(X_K, X_W \mid \text{do}(X_J))$  are the same. The former is obtained as:

$$P_M(X_K, X_W \mid \text{do}(X_J)) = (\text{pr}_{K \cup W} \circ (g, \text{id}_{\mathcal{X}_W}))_*(P) = (g_K, \text{id}_{\mathcal{X}_W})_*(P),$$

where  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  is the unique solution function of  $M$  and  $\text{pr}_{K \cup W} : \mathcal{X}_{V \cup W} \rightarrow \mathcal{X}_{K \cup W}$  is the canonical projection on the  $K \cup W$  components. The latter is obtained as:

$$P_{M_{\setminus L}}(X_K, X_W \mid \text{do}(X_J)) = (g_K, \text{id}_{\mathcal{X}_W})_*(P)$$

since by Lemma 6.7.2,  $g_K$  is the (unique) solution function of  $M_{\setminus L}$ . This means that both push-forwards are identical.

Let  $T \subseteq K \cup W$ . Then  $(M_{\setminus L})_{\text{do}(T)} = (M_{\text{do}(T)})_{\setminus L}$  by Proposition 6.7.4. The observable equivalence of  $M_{\text{do}(T)}$  and  $(M_{\text{do}(T)})_{\setminus L}$  w.r.t.  $(K \cup W) \setminus T$  hence implies the observable equivalence of  $M_{\text{do}(T)}$  and  $(M_{\setminus L})_{\text{do}(T)}$  w.r.t.  $(K \cup W) \setminus T$ . Since this holds for all  $T \subseteq K \cup W$ ,  $M$  and  $M_{\setminus L}$  are interventionally equivalent w.r.t.  $K \cup W$ .  $\square$

So the marginalization operation indeed effectively hides the details of a subsystem, while preserving the causal semantics on the remaining part.

## 6.8. Graphs of SCMs

A useful abstraction of an SCM is its graph. The directed edges in the graph of an SCM will express the following “parent”-relation that captures *functional dependencies* in the structural equations / causal mechanisms.

**Definition 6.8.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. For  $i \in J \cup V \cup W$  and  $j \in V$ , we say that  $i$  is a parent of  $j$  according to  $M$  if there does not exist a measurable function  $\tilde{f}_j : \mathcal{X}_{(J \cup V \cup W) \setminus \{i\}} \rightarrow \mathcal{X}_j$  such that for all  $x \in \mathcal{X}$ ,*

$$x_j = f_j(x) \iff x_j = \tilde{f}_j(x_{\setminus i}),$$

where  $x_{\setminus i}$  is shorthand for  $x_{(J \cup V \cup W) \setminus \{i\}}$ .

In words,  $i$  is parent of  $j$  if the causal mechanism for  $j$  is *not* equivalent to one that is constant with respect to its input component  $i$ . By definition, exogenous (input and random) variables have no parents. Note that the parent-relationship is preserved under equivalence. Using directed edges to encode the parent-relationship, we define the graph of the SCM.

**Definition 6.8.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. The CDG  $(J, V \cup W, E)$  with input nodes  $J$ , output nodes  $V \cup W$ , and directed edges*

$$E = \{i \rightarrow j : i \in J \cup W \cup V, j \in V : i \text{ is parent of } j \text{ according to } M\}$$

*is called the graph of the SCM and will be denoted as  $G(M)$ .*

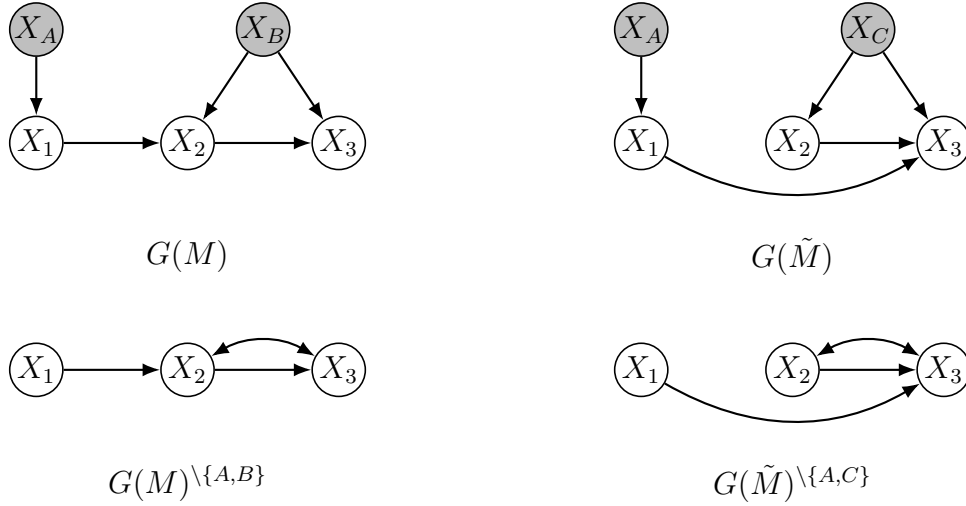


Figure 16: The graphs of the interventionally equivalent SCMs  $M$  (left) and  $\tilde{M}$  (center) corresponding to Example 6.8.3. If we consider only the endogenous variables  $X_1, X_2, X_3$  to be observed, and the other (exogenous random) variables to be latent, we can represent this with the two marginal graphs in the bottom row.

Since the parent-relationship is preserved under equivalence, equivalent SCMs have the same graphs. However, observably equivalent SCMs may have different graphs, and even interventionally equivalent SCMs may have different graphs.

**Example 6.8.3.** Consider the acyclic SCM  $M$  with endogenous variables  $X_1, X_2, X_3$  with co-domains  $\{-1, 1\}, \{-1, 1\}, \{-2, 0, 2\}$ , respectively, and structural equations

$$\begin{aligned} X_1 &= X_A \\ X_2 &= X_1 X_B \\ X_3 &= X_2 + X_B, \end{aligned}$$

with independent exogenous random variables  $X_A, X_B \sim \text{Uni}(\{-1, 1\})$ . Its graph  $G(M)$  is depicted in Figure 16 (left).

Consider also the acyclic SCM  $\tilde{M}$  with endogenous variables  $X_1, X_2, X_3$  with co-domains  $\{-1, 1\}, \{-1, 1\}, \{-2, 0, 2\}$ , respectively, and structural equations

$$\begin{aligned} X_1 &= X_A \\ X_2 &= X_C \\ X_3 &= X_2 + X_1 X_C, \end{aligned}$$

with independent exogenous random variables  $X_A, X_C \sim \text{Uni}(\{-1, 1\})$ . Its graph  $G(\tilde{M})$  is depicted in Figure 16 (center).  $\tilde{M}$  was obtained from  $M$  by making a change of variables  $x_C = x_1 x_B$ . One can check that  $\tilde{M}$  is interventionally equivalent to  $M$  with respect to  $\{1, 2, 3\}$  (that is, if we observe and can intervene on  $X_1, X_2, X_3$ , while the

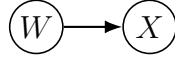
other variables are latent and non-intervenable), even though its the graph  $G(M)$  differs from the graph  $G(\tilde{M})$ , and the marginal graph  $G(M)^{\setminus\{A,B\}}$  differs from the marginal graph  $G(\tilde{M})^{\setminus\{A,C\}}$ .

Our goal was to allow for possible cycles in an SCM. This means that we may also encounter *self-cycles*.

**Example 6.8.4.** *The SCM in Example 6.6.4, with endogenous variable  $X$  and exogenous random variable  $W$  and structural equation*

$$X = -X^3 + X + W$$

has graph:



It does not have a self-cycle at  $X$ , since  $X$  is not a parent to itself: the structural equation is equivalent to

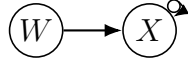
$$X = \sqrt[3]{W},$$

where  $X$  does not appear on the r.h.s..

On the other hand, the graph of the SCM in Example 6.6.5 with structural equation

$$X = WX + c$$

does have a self-cycle at  $X$ :



Self-cycles are a warning sign of complications with respect to solvability.

To understand why self-cycles are in some sense inevitable, consider an SCM with structural equations

$$X_1 = X_2$$

$$X_2 = X_3$$

$$X_3 = X_1$$

Marginalizing out  $X_2$  and  $X_3$  gives an SCM with structural equation

$$X_1 = X_1,$$

which turns the cycle  $X_1 \rightarrow X_3 \rightarrow X_2 \rightarrow X_1$  into a self-cycle  $X_1 \rightarrow X_1$ . Self-cycles are of no concern, though, when restricting to the class of simple SCMs.

**Proposition 6.8.5.** *For  $j \in V$ , we have that there is a self-cycle  $j \rightarrow j$  in  $G(M)$  if and only if  $M_{[j]}$  is not uniquely solvable. In particular, graphs of simple SCMs have no self-cycles.*

While the graph explicitly represents all exogenous random variables, coarser representations obtained via (graphical) marginalization are also useful.

**Remark 6.8.6.** An often used convention is to consider all exogenous random variables to be latent. In that case, a useful graph to consider is the marginalization  $G^{\setminus W}(M)$  of the graph  $G(M)$ . The marginal CDMG  $G^{\setminus W}(M) = (J, V, E, L)$  has input nodes  $J$ , output nodes  $V$ , directed edges

$$E = \{i \rightarrow j : i \in J \cup V, j \in V : i \text{ is parent of } j \text{ according to } M\}$$

and bidirected edges

$$L = \{j \leftrightarrow k : j \in V, k \in V, j \neq k : j \text{ and } k \text{ have common parent } i \in W \text{ according to } M\}.$$

In other words,  $G^{\setminus W}(M)$  is obtained from  $G(M)$  by replacing the nodes representing the exogenous random variables and their outgoing directed edges with bidirected edges, i.e., any pattern  $\leftarrow i \rightarrow$  with  $i \in W$  is replaced by  $\leftrightarrow$ .

We already provided definitions for hard interventions on graphs (Definition 3.2.1) and for marginalizations (latent projections) of graphs (Definition 3.2.14).

The mapping that maps an SCM to its graph is compatible with the elementary operations on SCMs.

**Proposition 6.8.7.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. Then

- Hard interventions: for  $T \subseteq J \cup V \cup W$ ,

$$G(M_{\text{do}(T)}) = G(M)_{\text{do}(T)}.$$

- Marginalizations: If  $M$  is simple, then for  $L \subseteq V$ ,

$$G(M_{\setminus L}) \subseteq G(M)^{\setminus L}.$$

*Proof.* The first statement follows by writing out the definitions. The second statement is somewhat more involved. We will first prove it in case  $L = \{\ell\}$  consists of a single node. Let  $G = G(M)$ . By using the definition of the parent relation (repeatedly), we can find a function  $\tilde{f}_\ell : \mathcal{X}_{\text{Pa}^G(\ell)} \rightarrow \mathcal{X}_\ell$  such that for all  $x \in \mathcal{X}$ :

$$x_\ell = f_\ell(x) \iff x_\ell = \tilde{f}_\ell(x_{\text{Pa}^G(\ell)}).$$

Since  $\ell \notin \text{Pa}^G(\ell)$  because  $M$  is simple, the unique solution function  $\tilde{g}_\ell : \mathcal{X}_{J \cup V \setminus \{\ell\} \cup W} \rightarrow \mathcal{X}_\ell$  of  $M_{\text{do}(V \setminus \{\ell\})}$  satisfies  $\tilde{g}_\ell(x_{\setminus \ell}) = \tilde{f}_\ell(x_{\text{Pa}^G(\ell)})$ , i.e., it only depends on the parents of  $\ell$ . When constructing the marginalized causal mechanism for  $M_{\setminus \{\ell\}}$ , we substitute  $\tilde{g}_\ell(x_{\setminus \ell})$  into the  $\ell$ 'th input of the causal mechanism  $f_j$  of  $M$ , for  $j \in V \setminus \{\ell\}$ . Since  $\tilde{g}_\ell$  only depends on  $\text{Pa}^G(\ell)$ , we get that  $\text{Pa}^{\tilde{G}}(j) \subseteq \text{Pa}^G(j) \setminus \{\ell\} \cup \text{Pa}^G(\ell)$ , where  $\tilde{G} = G(M_{\setminus \ell})$ . But we also have  $\text{Pa}^{G^{\setminus \{\ell\}}}(j) = \text{Pa}^G(j) \setminus \{\ell\} \cup \text{Pa}^G(\ell)$  by definition of the graphical marginalization. Hence  $\text{Pa}^{\tilde{G}}(j) \subseteq \text{Pa}^{G^{\setminus \{\ell\}}}$  for all  $j \in V \setminus \{\ell\}$ , and we have shown that  $G(M_{\setminus \{\ell\}}) \subseteq G(M)^{\setminus \{\ell\}}$ . For the general case, we can make use of induction and the facts that both for graphs and simple SCMs, we can obtain a marginalization over a subset by repeatedly marginalizing out a single remaining node in the subset, in arbitrary order.  $\square$

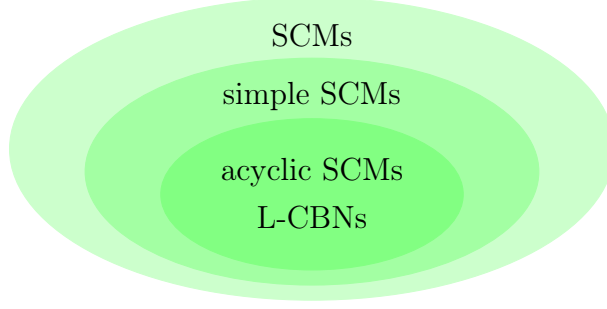


Figure 17: Venn diagram for different causal modeling classes.

A subclass of SCMs that is often considered are acyclic SCMs.

**Definition 6.8.8.** *An SCM  $M$  is called acyclic if its graph  $G(M)$  is acyclic.*

If one models static systems, then using acyclic SCMs rules out the presence of causal cycles (e.g., feedback loops) in the system. Acyclic SCMs are a subclass of the more general class of simple SCMs.

**Proposition 6.8.9.** *Acyclic SCMs are simple.*

*Proof.* We first show that acyclic SCMs are uniquely solvable. Let  $M$  be an acyclic SCM. Its graph  $G := G(M)$  is acyclic, and hence has a topological order  $<$ . Consider  $f_v$ , the causal mechanism for  $v \in V$ . The parents  $\text{Pa}^G(v)$  precede  $v$  in the topological order. Since  $f_v$  can be rewritten to be constant in the non-parents of  $v$  (similar to how this was done in the proof of Proposition 6.8.7), we can consider  $f_v : \mathcal{X} \rightarrow \mathcal{X}_v$  as a function  $f_v : \mathcal{X}_{\text{Pred}_{<}^G(v)} \rightarrow \mathcal{X}_v$  instead. We can then inductively define the components

$$g_v : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_v : (x_J, x_W) \mapsto f_v(g_{\text{Pred}_{<}^G(v)}(x_J, x_W))$$

that together form a solution function  $g : \mathcal{X}_{J \cup W} \rightarrow \mathcal{X}_V$ . This construction also exhibits the uniqueness of  $g : \mathcal{X}_{J \cup W} \rightarrow \mathcal{X}_V$ .

Next consider  $M_{\text{do}(T)}$ , the intervened SCM for a hard intervention on  $M$  with target  $T \subseteq V \cup W \cup J$ . It has graph  $G(M_{\text{do}(T)}) = G_{\text{do}(T)}$ , whose edges form a subset of the edges of  $G$  (but where some output nodes have become input nodes), and hence is also acyclic. Therefore, also  $M_{\text{do}(T)}$  is uniquely solvable. Since this holds for all targets  $T \subseteq V \cup W \cup J$ , we conclude that  $M$  is simple.  $\square$

## 6.9. Interventional equivalence of acyclic SCMs and (L-)CBNs

Figure 17 shows a Venn diagram to illustrate the relationships between the different classes of causal models that we introduced. We will show that in a precise sense, acyclic SCMs and L-CBNs are equally expressive.

**Definition 6.9.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM and let  $\tilde{M}$  be an L-CBN with graph  $\tilde{G} = (\tilde{J}, \tilde{V} \cup \tilde{U}, \tilde{E})$ , spaces  $\tilde{\mathcal{X}}_{\tilde{v}}$  for  $\tilde{v} \in \tilde{V} \cup \tilde{U} \cup \tilde{J}$ , and Markov kernels*

$$P_{\tilde{v}}(X_{\tilde{v}} | X_{\text{Pa}_{\tilde{G}}(\tilde{v})}).$$



We say that  $M$  is interventionally equivalent to  $\tilde{M}$  w.r.t.  $O \subseteq (V \cup W) \cap \tilde{V}$  if  $\mathcal{X}_O = \tilde{\mathcal{X}}_O$ ,  $\mathcal{X}_{J \cap \tilde{J}} = \tilde{\mathcal{X}}_{J \cap \tilde{J}}$ , and for any hard intervention  $\text{do}(T)$  with  $T \subseteq O$ , the intervened Markov kernel  $P_M(X_{O \setminus T} | \text{do}(X_{J \cup T}))$  of  $M$  equals the intervened Markov kernel  $P(X_{O \setminus T} | \text{do}(X_{\tilde{J} \cup T}))$  of  $\tilde{M}$ .

**Proposition 6.9.2.** *i) Given an acyclic SCM  $M = (J, V, W, \mathcal{X}, P, f)$  and an observed subset  $O \subseteq V \cup W$ , we can construct an L-CBN  $\tilde{M}$  with observed output variables  $\tilde{V} = O$ , latent output variables  $\tilde{U} = (V \cup W) \setminus O$ , input variables  $J$ , graph  $G(M)$ , and spaces  $\mathcal{X}_v$  for  $v \in V \cup W \cup J$  that is interventionally equivalent to  $M$  w.r.t.  $O$ .*

*ii) Given an L-CBN  $\tilde{M} = \left( G^+ = (J, (O, U), E^+), \left( P_v(X_v | X_{\text{Pa}^{G^+}(v)}) \right)_{v \in O \cup U} \right)$  with observed variables  $O$ , we can construct an acyclic SCM  $M$  with input variables  $J$ , endogenous variables  $O \dot{\cup} U$  and marginal graph  $G^{O \cup U}(M) = G^+$  that is interventionally equivalent to  $\tilde{M}$  w.r.t.  $O$ .*

*Proof.* i) Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an acyclic SCM with graph  $G(M) = (J, V \cup W, E)$ . Define  $\tilde{M}$  as the L-CBN with observed output variables  $\tilde{V} = O$ , latent output variables  $\tilde{U} = (V \cup W) \setminus O$ , input variables  $J$ , graph  $G^+ := G(M)$ , spaces  $\mathcal{X}_v$  for  $v \in V \cup W \cup J$ , and the following Markov kernels. For  $v \in V$ , we write its structural equation as:

$$x_v = f_v(x_{\text{Pa}^{G^+}(v)})$$

with  $f_v : \mathcal{X}_{\text{Pa}^{G^+}(v)} \rightarrow \mathcal{X}_v$ , where  $v \notin \text{Pa}^{G^+}(v)$  because the graph is acyclic. We then define the corresponding (deterministic) Markov kernel

$$P_v \left( X_v | X_{\text{Pa}^{G^+}(v)} \right) := \delta_{f_v}(X_v | X_{\text{Pa}^{G^+}(v)}),$$

encoding the causal mechanisms of  $M$ . The Markov kernels for  $w \in W$  are defined as:

$$P_w \left( X_w | X_{\text{Pa}^{G^+}(w)} \right) := P_w(X_w),$$

encoding the exogenous distributions of  $M$ , where we note that  $\text{Pa}^{G^+}(w) = \emptyset$ . One can check that this L-CBN  $\tilde{M}$  does the job.

ii) Let

$$\tilde{M} = \left( G^+ = (J, (O, U), E^+), \left( P_v(X_v | X_{\text{Pa}^{G^+}(v)}) \right)_{v \in O \cup U} \right)$$

be an L-CBN. For every  $v \in O \dot{\cup} U$ , we can write the Markov kernel  $P_v$  as the composition of a deterministic one and a uniform distribution  $P_{\bar{v}}(X_{\bar{v}})$  on  $\mathcal{X}_{\bar{v}} := [0, 1]$  by Remark 2.7.4:

$$P_v(X_v | X_{\text{Pa}^{G^+}(v)}) = \delta(f_v | X_{\bar{v}}, X_{\text{Pa}^{G^+}(v)}) \circ P_{\bar{v}}(X_{\bar{v}})$$

for some measurable function  $f_v : \mathcal{X}_{\bar{v}} \times \mathcal{X}_{\text{Pa}^{G^+}(v)} \rightarrow \mathcal{X}_v$ . Here, we introduced new variables  $\bar{v}$  for each  $v \in O \dot{\cup} U$ .

Define now the SCM  $M = (J, O \dot{\cup} U, \bar{O} \dot{\cup} \bar{U}, \mathcal{X}_J \times \mathcal{X}_{O \dot{\cup} U} \times \mathcal{X}_{\bar{O} \dot{\cup} \bar{U}}, P, f)$  by taking

$$P(X_{\bar{O} \dot{\cup} \bar{U}}) = \bigotimes_{\bar{v} \in \bar{O} \dot{\cup} \bar{U}} P_{\bar{v}}(X_{\bar{v}})$$

and

$$f = (f_v)_{v \in O \dot{\cup} U}.$$

That is, the uniformly distributed random variables  $X_{\bar{v}}$  become the exogenous random variables, all independent and uniformly distributed on  $[0, 1]$ , and the components of the causal mechanism correspond to the deterministic functions used to represent the Markov kernels. This SCM does the job, as one can check.  $\square$

Simple SCMs are more expressive than acyclic SCMs because they can model (sufficiently weak) causal cycles. SCMs in general are even more expressive because they can also model stronger cycles that not necessarily lead to unique solvability under any hard intervention, but this generality comes with a substantially increased complexity of the theory and interpretability. Simple SCMs form a “sweet spot” in the sense that they allow cyclic relationships yet their theory is not much more complicated than that of acyclic SCMs: the main difference consists in replacing  $d$ -separation with  $\sigma$ -separation.

Even SCMs may not be the ultimate way of modeling cyclic causal systems. Indeed, for such systems, it might be that the conceptual notion of interventions *targeting variables* is misguided in general, and perhaps should be replaced by the notion of intervening on *functional constraints* [BvDM21].

## 6.10. Examples

In many systems occurring in the real world, feedback loops between observed variables are present. Such systems can often be described by a system of (random) differential equations. The equilibrium states of such systems can sometimes be causally modelled by an SCM [BM18].

For illustration purposes we provide two examples, the first consisting of interacting masses that are attached to springs that can be described at equilibrium with a simple SCM, the second being the famous price-supply-demand model that has been very popular in econometrics, and which corresponds to a non-simple SCM at equilibrium.

**Example 6.10.1** (Damped coupled harmonic oscillator). *Consider a one-dimensional system of  $d$  masses  $m_i \in \mathbb{R}$  ( $i = 1, \dots, d$ ) with positions  $Q_i$ . The masses are coupled by springs, with spring constants  $k_i > 0$  ( $i = 0, \dots, d$ ) and equilibrium lengths  $\ell_i > 0$  ( $i = 0, \dots, d - 1$ ), under influence of friction with friction coefficients  $b_i > 0$  ( $i = 1, \dots, d$ ). The endpoints are considered fixed at positions  $Q_0 < Q_{d+1}$  (see Figure 18 (top)). From elementary physics, we know that the equations of motion of this system are given by the following differential equations*

$$\frac{d^2 Q_i}{dt^2} = \frac{k_i}{m_i} (Q_{i+1} - Q_i - \ell_i) + \frac{k_{i-1}}{m_i} (Q_{i-1} - Q_i + \ell_{i-1}) - \frac{b_i}{m_i} \frac{dQ_i}{dt} \quad i = 1, \dots, d.$$

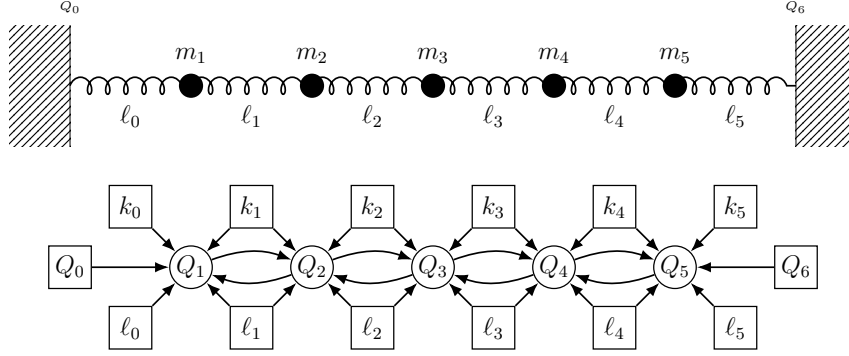


Figure 18: Damped coupled harmonic oscillator (top) and the graph of the SCM that describes the positions of the masses at equilibrium (bottom) of Example 6.10.1 for  $d = 5$ , where the spring lengths and constants are considered as exogenous input variables.

The dynamics of the masses, in terms of the position  $Q_i$ , velocity  $\frac{dQ_i}{dt}$  and acceleration  $\frac{d^2Q_i}{dt^2}$ , is described by a single and separate equation of motion for each mass. Under friction, i.e.,  $b_i > 0$  ( $i = 1, \dots, d$ ), there is a unique equilibrium position, where the sum of forces vanishes for each mass. If one moves one or several masses out of their equilibrium positions and releases them, then the masses will start to oscillate, but eventually these oscillations dampen out and the masses converge to their unique equilibrium position. At equilibrium (i.e., for  $t \rightarrow \infty$ ) the velocity  $\frac{dQ_i}{dt}$  and acceleration  $\frac{d^2Q_i}{dt^2}$  of the masses vanish (i.e.,  $\frac{dQ_i}{dt}, \frac{d^2Q_i}{dt^2} \rightarrow 0$ ), and thus the following equation holds at equilibrium

$$0 = \frac{k_i}{m_i}(Q_{i+1} - Q_i - \ell_i) + \frac{k_{i-1}}{m_i}(Q_{i-1} - Q_i + \ell_{i-1})$$

for each mass ( $i = 1, \dots, d$ ). By solving each of these equations w.r.t.  $Q_i$ , we obtain that the equilibrium positions  $Q_i$  of the masses are given by

$$Q_i = \frac{k_i(Q_{i+1} - \ell_i) + k_{i-1}(Q_{i-1} + \ell_{i-1})}{k_i + k_{i-1}}.$$

By considering the  $\ell_i$ ,  $k_i$  and  $Q_0$  and  $Q_{d+1}$  as exogenous (input or random) variables, and the  $Q_i$  ( $i = 1, \dots, d$ ) as endogenous variables, we arrive at an SCM with causal mechanism

$$f_i(q, \ell, k) = \frac{k_i(q_{i+1} - \ell_i) + k_{i-1}(q_{i-1} + \ell_{i-1})}{k_i + k_{i-1}}.$$

for  $i = 1, \dots, d$ . Its graph is depicted in Figure 18 (bottom). This SCM allows us to describe the equilibrium behavior of the system under perfect intervention. For example, when forcing the mass  $j$  to a fixed position  $Q_j = \xi_j$  with  $0 \leq \xi_j \leq L$ , the equilibrium positions of the masses correspond to the solutions of the intervened model  $M_{\text{do}(Q_j=\xi_j)}$ .

**Exercise 6.10.2.** Prove that the SCM that describes the equilibrium states of a damped coupled harmonic oscillator is simple (see also Proposition 6.5.4). Hint: you can use that

the determinant of a tridiagonal matrix of the following form is given by the expression on the r.h.s.:

$$\det \begin{pmatrix} k_0 + k_1 & -k_1 & & & & \\ -k_1 & k_1 + k_2 & -k_2 & & & \\ & -k_2 & k_2 + k_3 & \ddots & & \\ & & \ddots & \ddots & -k_{d-1} & \\ & & & -k_{d-1} & k_{d-1} + k_d & \end{pmatrix} = \sum_{i=0}^d \prod_{\substack{j=0 \\ j \neq i}}^d k_j$$

Next, we show that the well-known market equilibrium model from economics, can be described by a (non-simple) SCM. This example illustrates how self-cycles enrich the class of SCMs.

**Example 6.10.3** (Price, supply and demand). *Let  $D$  denote the demand and  $S$  the supply of a quantity of a product. The price of the product is denoted by  $R$ . The following system of differential equations describes how the demanded and supplied quantities are determined by the price, and how price adjustments occur in the market:*

$$\begin{aligned} D &= \beta_D R + E_D \\ S &= \beta_S R + E_S \\ \frac{dR}{dt} &= D - S, \end{aligned}$$

where  $E_D$  and  $E_S$  are exogenous random influences on the demand and supply respectively,  $\beta_D < 0$  is the reciprocal of the slope of the demand curve, and  $\beta_S > 0$  is the reciprocal of the slope of the supply curve. At the situation known as a “market equilibrium”, the price is determined implicitly by the condition that demanded and supplied quantities should be equal, since  $\frac{dR}{dt} = 0$  at equilibrium. At equilibrium, hence, we obtain an SCM  $M$  with causal mechanism defined by:

$$\begin{aligned} f_D(d, s, r, e_D, e_S) &:= \beta_D r + e_D \\ f_S(d, s, r, e_D, e_S) &:= \beta_S r + e_S \\ f_R(d, s, r, e_D, e_S) &:= r + (d - s). \end{aligned}$$

Note how we use a self-cycle for  $r$  in order to implement the equilibrium equation  $d = s$  as the causal mechanism for the price  $r$ . Its graph is depicted in Figure 19 (left).

**Exercise 6.10.4.** *Prove that the SCM  $M$  that describes the equilibrium states of the price-supply-demand model is uniquely solvable, but not simple. Consider the following interventions:  $\text{do}(D = \delta)$ ,  $\text{do}(S = \sigma)$ ,  $\text{do}(R = \rho)$ , and all possible combinations thereof. Which of (the combinations of) these interventions give an intervened SCM that is still uniquely solvable? Which of these interventions on the SCM correspond with the equilibrium state of a similarly intervened market dynamics model? Summarizing: could this be a realistic causal equilibrium model of an ideal market, or is there something wrong with it (perhaps due to the self-cycle)?*

*(Bonus: can you model the market equilibrium with an SCM without self-cycles?)*

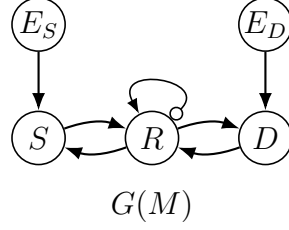


Figure 19: The graph of the SCM  $M$  of Example 6.10.3.

While the price-supply-demand example shows that not all cyclic SCMs that occur “in the wild” are simple, we have chosen to restrict ourselves mostly to simple SCMs for this lecture. Generalizations of the theory presented here for simple SCMs to non-simple ones are provided in [BFPM21].

We will finish this section by showing that SCMs are simple if the causal mechanisms are sufficiently weak and smooth.

**Proposition 6.10.5.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM with real-valued endogenous variables, that is,  $\mathcal{X}_v = \mathbb{R}$  for each  $v \in V$ . If for each subset  $U \subseteq V$ , and for all values  $x_W \in \mathcal{X}_W$ ,  $x_J \in \mathcal{X}_J$ ,  $x_{V \setminus U} \in \mathcal{X}_{V \setminus U}$ , the mapping*

$$\mathcal{X}_U \rightarrow \mathcal{X}_U : x_U \mapsto f(x_J, x_{V \setminus U}, x_U, x_W)$$

*is Lipschitz continuous with Lipschitz constant  $L_U(x_J, x_{V \setminus U}, x_W) < 1$  with respect to some norm  $\|\cdot\|$ , then  $M$  is simple.*

*Proof.* By definition,  $M$  is simple if for all  $U \subseteq V$ , for all  $x_W \in \mathcal{X}_W$ , for all  $x_J \in \mathcal{X}_J$ , for all  $x_{V \setminus U} \in \mathcal{X}_{V \setminus U}$ , the equation

$$x_U = f(x_J, x_{V \setminus U}, x_U, x_W) \tag{42}$$

has a unique solution for  $x_U \in \mathcal{X}_U$ . This is a fixed point equation for  $x_U$ , and hence it has a unique solution by Banach’s fixed point theorem if it is a contraction (Lipschitz continuous with Lipschitz constant  $< 1$  with respect to some norm  $\|\cdot\|$ ).  $\square$

**Remark 6.10.6.** *This also provides us with a method for sampling from a simple SCM that satisfies the assumption in Proposition 6.10.5. The solution to equation (42) can be obtained by iterating the updates*

$$x_U^{(n+1)} = f(x_J, x_{V \setminus U}, x_U^{(n)}, x_W)$$

*until convergence.*

To give a more concrete example: a class of SCMs that satisfies the contractivity condition is given by neural networks with sufficiently weak weights.

**Example 6.10.7.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM with real-valued variables, that is,  $\mathcal{X}_k = \mathbb{R}$  for each  $k \in V \cup W \cup J$ . Suppose that the causal mechanism is of the form

$$f_u = h \left( \sum_{j \in J} A_{uj} x_j + \sum_{w \in W} A_{uw} x_w + \sum_{v \in V} A_{uv} x_v + b_u \right), \quad u \in V,$$

with weights  $A \in \mathbb{R}^{V \times (V \cup W \cup J)}$ , biases  $b \in \mathbb{R}^V$  and activation function  $h : \mathbb{R} \rightarrow \mathbb{R}$ .

The conditions in Proposition 6.10.5 are satisfied if the following conditions both hold:

1.  $\sup_{x \in \mathbb{R}} |h'(x)| \leq C$  with  $0 < C < \infty$ , and
2.  $\|A_{UU}\| < \frac{1}{C}$  for every subset  $U \subseteq V$  of cardinality  $\#(U) \geq 2$ , where  $\|\cdot\|$  can be one of the matrix norms:  $\|\cdot\|_p$ ,  $p \geq 1$ , or  $\|\cdot\|_\infty$ .

*Proof.* By the mean value theorem, it suffices to show that for every subset  $U \subseteq V$  of cardinality  $\#(U) \geq 2$  and every value  $(x_J, x_W, x_{V \setminus U})$  the partial derivative is bounded:

$$\sup_{x_U \in \mathcal{X}_U} \left\| \frac{\partial f_U}{\partial x_U}(x_J, x_{V \setminus U}, x_U, x_W) \right\| \leq L_U(x_J, x_{V \setminus U}, x_W) < 1$$

for  $\|\cdot\|$  some matrix norm. In our case we have:

$$\frac{\partial f_U}{\partial x_U}(x_J, x_{V \setminus U}, x_U, x_W) = \text{diag}(\eta)_{UU} A_{UU}.$$

where  $\eta$  is a vector in  $\mathbb{R}^V$  with entries

$$\eta_v = h'(A_v(x_J, x_V, x_W)^\top + b_v).$$

If  $|h'(x)| \leq C < \infty$  for all  $x \in \mathbb{R}$ , and  $\|\cdot\|$  is either  $\|\cdot\|_p$ ,  $p \geq 1$ , or  $\|\cdot\|_\infty$ , then  $\|\text{diag}(\eta)_{UU}\| \leq C$ . Since  $\|A_{UU}\| < \frac{1}{C}$  we get

$$\|\text{diag}(\eta)_{UU} A_{UU}\| \leq \|\text{diag}(\eta)_{UU}\| \cdot \|A_{UU}\| =: L_U(x_J, x_{V \setminus U}, x_W) < C \frac{1}{C} = 1.$$

□

**Remark 6.10.8.** Note that we can put  $C = 1$  for popular activation functions  $h(x)$  like  $\tanh(x)$ ,  $\text{ReLU}(x) = \max(0, x)$ ,  $\sigma(x) = \frac{1}{1 + \exp(-x)}$ , *LeakyRelu*, and *SoftPlus* ( $x) = \ln(1 + e^x)$ ). Further note that by using one of these activation functions  $h(x)$  and  $\|\cdot\| = \|\cdot\|_\infty$  all the conditions are satisfied if we choose the weights  $A_{v,k}$  such that for all  $v \in V$ :

$$\sum_{k \in V \cup W \cup J} |A_{v,k}| < 1,$$

and  $A_{v,v} = 0$ . While this is far from necessary, it is easy to check.

## 7. Markov property for simple SCMs

By making use of the ‘acyclification’, we can extend the global Markov property for L-CBNs to a global Markov property for simple SCMs. The difference is that the Markov property for SCMs is formulated in terms of  $\sigma$ -separation rather than  $d$ -separation. With the help of this Markov property, we can derive a very analogous theory for simple SCMs as for L-CBNs, with a do-calculus and adjustment.

### 7.1. Acyclifications

In Section 3.5, we defined acyclifications of a CDMG. We can also define an operation with the same name on SCMs.

**Definition 7.1.1** (Acyclification of SCM). *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G(M)$ . For each strongly connected component  $C$  of  $G(M)$  (i.e., a set of the form  $\text{Sc}^{G(M)}(v)$  for  $v \in V$ ), let  $g_C : \mathcal{X}_{J \cup (V \setminus C) \cup W} \rightarrow \mathcal{X}_C$  be the unique solution function for  $M_{[C]}$ . Define  $\tilde{f} : \mathcal{X}_J \times \mathcal{X}_V \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  by its components*

$$\tilde{f}_v(x_J, x_V, x_W) = (g_{\text{Sc}^{G(M)}(v)})_v(x_J, x_{V \setminus C}, x_W)$$

for  $v \in V$ .  $M^{\text{acy}} = (J, V, W, \mathcal{X}, P, \tilde{f})$  is called the acyclification of  $M$ .

The crucial property of this definition is the following result, which also motivates its name.

**Proposition 7.1.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Its acyclification  $M^{\text{acy}}$  is acyclic and observably equivalent to  $M$ .*

*Proof.* We construct a directed graph  $S$  from  $G := G(M)$  with its strongly connected components  $\{\text{Sc}^G(v) : v \in V \cup J \cup W\}$  as nodes, and directed edges  $C \rightarrow D$  if there is a directed edge  $c \rightarrow d$  in  $G$  with  $c \in C, d \in D$  and  $C \neq D$ . The graph  $S$  cannot contain a directed cycle, as that would imply the existence of a directed cycle in  $G$  that traverses more than one of its strongly connected components. Hence  $S$  is a DAG.

Choose a topological ordering  $<$  of  $S$ . Any node  $C$  in  $S$  can only have incoming directed edges in  $S$  from  $\text{Pred}_{<}^S(C)$ . This implies that for  $v \in V$ ,  $C = \text{Sc}^G(v)$  can only have incoming edges in  $G$  from  $\bigcup \text{Pred}_{<}^S(C)$ . That implies that the causal mechanism  $f_C$  can only depend on variables in  $\bigcup \text{Pred}_{<}^S(C)$ , and hence the unique solution function  $g_C$ , and therefore  $\tilde{f}_C$ , can depend on variables in  $\bigcup \text{Pred}_{<}^S(C)$  only. Therefore, for  $v, w \in V$ , a directed edge  $w \rightarrow v$  in  $G(M^{\text{acy}})$  implies  $w \in \bigcup \text{Pred}_{<}^S(\text{Sc}^G(v))$ . We can therefore refine the topological ordering  $<$  of  $S$  to a topological ordering of  $G(M^{\text{acy}})$ , by arbitrarily ordering the nodes within each strongly connected component of  $G$ . Hence  $G(M^{\text{acy}})$  is acyclic.

$M$  and  $M^{\text{acy}}$  are observably equivalent by construction: for all  $x \in \mathcal{X}$ ,

$$\begin{aligned}
x &= \tilde{f}(x) \\
&\iff \forall C \in S \cap V : x_C = \tilde{f}_C(x) \\
&\iff \forall C \in S \cap V : x_C = g_C(x_J, x_{V \setminus C}, x_W) \\
&\iff \forall C \in S \cap V : x_C = f_C(x) \\
&\iff x = f(x).
\end{aligned}$$

□

The notion of acyclification of an SCM is compatible with that of a graph:

**Proposition 7.1.3.** *Let  $M$  be a simple SCM. Then  $G(M^{\text{acy}}) \subseteq G'$  for any acyclification  $G'$  of  $G(M)$ .*

*Proof.* By definition, the two graphs have the same nodes (input nodes  $J$  and output nodes  $V \cup W$ ).  $G(M^{\text{acy}})$  has no bidirected edges, but  $G'$  might. If there is a directed edge  $i \rightarrow j$  in  $G(M^{\text{acy}})$  with  $i \in J \cup V \cup W$  and  $j \in V$ , then the solution function  $g_{\text{Sc}^G(j)}$  of  $M_{[\text{Sc}^G(j)]}$  depends on  $x_i$ . This can only happen if  $i \notin \text{Sc}^G(j)$  and  $i$  is a parent of some  $k$  according to  $M$  with  $k \in \text{Sc}^G(j)$ , i.e., if  $i \rightarrow k$  in  $G(M)$ . In that case,  $i \rightarrow j$  in  $G'$  by definition of the graphical acyclification. □

Hence, two nodes in the same strongly connected component of  $G(M)$  do not have any edge between them in  $G(M^{\text{acy}})$ , whereas they necessarily have a connecting edge in any acyclification  $G'$  of  $G(M)$ . For two nodes in different strongly connected components of  $G(M)$ , the edges in  $G(M^{\text{acy}})$  are also present in  $G'$ , but not necessarily vice versa, as some parent-relations may cancel out in the solution function.

## 7.2. Global Markov property for simple SCMs

With the help of the acyclifications, we can easily derive a Markov property for simple SCMs from the Markov property for CBNs by reducing the general cyclic case to an acyclic case.

**Corollary 7.2.1** (Global Markov property for simple SCMs). *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G(M)$  and Markov kernel  $P_M(X_V, X_W \mid \text{do}(X_J))$ . Then for all  $A, B, C \subseteq J \cup V \cup W$  (not necessarily disjoint) we have the implication:*

$$A \underset{G(M)}{\perp}^{\sigma} B \mid C \implies X_A \underset{P_M(X_V, X_W \mid \text{do}(X_J))}{\perp\!\!\!\perp} X_B \mid X_C.$$

*If one wants to make the implicit dependence on  $J$  more explicit one can equivalently also write:*

$$A \underset{G(M)}{\perp}^{\sigma} J \cup B \mid C \implies X_A \underset{P_M(X_V, X_W \mid \text{do}(X_J))}{\perp\!\!\!\perp} X_J, X_B \mid X_C.$$



*Proof.* Choose an acyclification  $G'$  of  $G(M)$ . Then:

$$\begin{aligned}
A \underset{G(M)}{\perp}^{\sigma} B | C &\iff A \underset{G'}{\perp}^d B | C \\
&\implies A \underset{G(M^{\text{acy}})}{\perp}^d B | C \\
&\implies X_A \underset{P_{M^{\text{acy}}}(X_V, X_W | \text{do}(X_J))}{\perp\!\!\!\perp} X_B | X_C \\
&\iff X_A \underset{P_M(X_V, X_W | \text{do}(X_J))}{\perp\!\!\!\perp} X_B | X_C.
\end{aligned}$$

For the various implications / equivalences, we used:

1.  $G'$  is an acyclification of  $G(M)$  together with Proposition 3.5.2;
2.  $G(M^{\text{acy}}) \subseteq G'$  from Proposition 7.1.3, and that removing edges cannot turn a  $d$ -separation into a  $d$ -connection;
3. the global Markov property Theorem 4.2.1 for  $M^{\text{acy}}$  interpreted as a causal Bayesian network as in the proof of Proposition 6.9.2 point i) (with deterministic Markov kernels for the endogenous variables, and purely probabilistic Markov kernels for the exogenous random variables), exploiting Proposition 7.1.2 that states that the acyclification  $M^{\text{acy}}$  is acyclic;
4. by Proposition 7.1.2, the acyclification  $M^{\text{acy}}$  has the same Markov kernel as the original SCM  $M$ .

□

### 7.3. Do-calculus for simple SCMs

With the global Markov property for simple SCMs, it becomes straightforward to derive the do-calculus for simple SCMs. First we will introduce some notation. The setting will be that a simple SCM  $M = (J, V, W, \mathcal{X}, P, f)$  is given. As an auxiliary way to model hard interventions  $\text{do}(B)$  for  $B \subseteq V \cup W$  that leads to easy rules in the do-calculus, we introduce additional intervention variables  $I_b$  for  $b \in B$ . We denote  $I_B := (I_b)_{b \in B}$ . This gives an extended SCM  $\tilde{M} = (J \cup I_B, V, W, \mathcal{X} \times \prod_{b \in B} (\mathcal{X}_b \cup \{\star\}), P, \tilde{f})$  with causal mechanism with components

$$\tilde{f}_b(x_{V \cup J \cup W}, x_{I_B}) := \begin{cases} f_b(x_{V \cup J \cup W}) & x_{I_b} = \star \\ x_{I_b} & x_{I_b} \in \mathcal{X}_b \end{cases}$$

for  $b \in B \cap V$ , and  $\tilde{f}_v(x) := f_v(x_{V \cup J \cup W})$  for  $v \in V \setminus B$ . Here,  $x_{I_b} = \star$  encodes that there is no intervention on  $b$ , while  $x_{I_b} \neq \star$  encodes that the hard intervention  $\text{do}(X_b = x_{I_b})$  is performed. We will denote the graph of  $M$  by  $G$  and the graph of  $\tilde{M}$  by  $G_{\text{do}(I_B)}$ . In the do-calculus, we make use of the extended graph  $G_{\text{do}(I_B)}$  to test the separation statement,

while the conclusion about the properties of certain Markov kernels concerns those of the original SCM  $M$ .

The Markov kernel of simple SCM  $M$  exists, is unique, and is denoted by  $P_M(X_{V \cup W} \mid \text{do}(X_J))$ . For a subset  $T \subseteq V \cup W$ , we write

$$P_M(X_{(V \cup W) \setminus T} \mid \text{do}(X_J, X_T)) := P_{M_{\text{do}(T)}}(X_{(V \cup W) \setminus T} \mid \text{do}(X_J, X_T)).$$

By conditioning on a subset  $S \subseteq (V \cup W) \setminus T$ , we obtain the conditional Markov kernel

$$P_M(X_{(V \cup W) \setminus (T \cup S)} \mid X_S, \text{do}(X_J, X_T)).$$

The only modification to the do-calculus for simple SCMs as compared to that for causal Bayesian networks is that we have to replace all  $d$ -separations by  $\sigma$ -separations.

**Theorem 7.3.1** (Almost-sure do-calculus for simple SCMs, simplified). *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G = G(M)$ . Assume that we have  $\sigma$ -finite reference measures  $\mu_v$  on  $\mathcal{X}_v$  for every  $v \in V \cup W$  and put  $\mu_F := \bigotimes_{v \in F} \mu_v$  for  $F \subseteq V \cup W$ . Let  $A, B, C \subseteq V \cup W$  and  $D \subseteq V \cup W \cup J$  be such that  $A, B, C, D$  are pairwise disjoint. Then we have the following 3 rules relating Markov kernels that can be generated from the SCM:*

1. Insertion/deletion of observation, for  $J \subseteq D$ : if

$$A \underset{G_{\text{do}(D)}}{\perp}^{\sigma} B \mid C \cup D, \quad \mu_{B \cup C} \ll P_M(X_B, X_C \mid \text{do}(X_D)) \ll \mu_{B \cup C},$$

then:

$$P_M(X_A \mid X_B, X_C, \text{do}(X_D)) = P_M(X_A \mid X_C, \text{do}(X_D)) \quad \mu_{B \cup C}\text{-a.s.}$$

2. Action/observation exchange, for  $J \subseteq D$ : if

$$A \underset{G_{\text{do}(I_B, D)}}{\perp}^{\sigma} I_B \mid B \cup C \cup D, \quad \mu_{B \cup C} \ll P_M(X_B, X_C \mid \text{do}(X_D)) \ll \mu_{B \cup C},$$

$$\mu_C \ll P_M(X_C \mid \text{do}(X_B, X_D)) \ll \mu_C,$$

then:

$$P_M(X_A \mid X_B, X_C, \text{do}(X_D)) = P_M(X_A \mid \text{do}(X_B), X_C, \text{do}(X_D)) \quad \mu_{B \cup C}\text{-a.s.}$$

3. Insertion/deletion of action, for  $J \subseteq D$ : if

$$A \underset{G_{\text{do}(I_B, D)}}{\perp}^{\sigma} I_B \mid C \cup D, \quad \mu_C \ll P_M(X_C \mid \text{do}(X_B, X_D)) \ll \mu_C,$$

$$\mu_C \ll P_M(X_C \mid \text{do}(X_D)) \ll \mu_C,$$

then

$$P_M(X_A \mid \text{do}(X_B), X_C, \text{do}(X_D)) = P_M(X_A \mid X_C, \text{do}(X_D)) \quad \mu_C\text{-a.s.}$$

4. Deletion of input: *If*

$$A \underset{G_{\text{do}(X_D)}}{\perp}^{\sigma} J \mid C \cup D, \quad \mu_C \ll P_M(X_C \mid \text{do}(X_{D \cup J})) \ll \mu_C,$$

*then there exists a Markov kernel  $P_M(X_A \mid X_C, \text{do}(X_D, \cancel{X_{J \setminus D}}))$  such that:*

$$P_M(X_A \mid X_C, \text{do}(X_D, \cancel{X_{J \setminus D}})) = P_M(X_A \mid X_C, \text{do}(X_{D \cup J})) \quad \mu_C\text{-a.s.}$$

*Proof.* The proof is analogous to that of Corollary 5.1.3, except that it applies the global Markov property for simple SCMs, Corollary 7.2.1, instead of the one for causal Bayesian networks, Theorem 4.2.1.  $\square$

While the derivation of the do-calculus relies essentially on the global Markov property, sometimes one can make use of the global Markov property and a more careful analysis of null sets to obtain stronger conclusions. In particular, Proposition 5.1.7 and Theorem 5.1.2 also hold for simple SCMs if one replaces the d-separation statements by the analogous  $\sigma$ -separation statements.

## 7.4. Adjustment

Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let us assume  $J \subseteq D$ . We are interested in estimating the conditional causal effect:

$$P_M(X_A \mid X_C, \text{do}(X_B, X_D)),$$

but we only have data from:

$$P_M(X_A, X_B, X_F \mid X_C, \text{do}(X_D)).$$

The following (pairwise disjoint) index sets will have the following roles:

$A \subseteq V$  : the outcome variables of interest.

$B \subseteq V \cup W$  : the treatment or intervention variables.

$C \subseteq V \cup W$  : general conditional (context) variables under which the data was collected.

$J \subseteq D \subseteq V \cup W \cup J$  : general interventional (context) variables that were set by the experimenter.

$F_0 \subseteq V \cup W$  : core adjustment variables, i.e. features that were measured.

$F_1 \subseteq V \cup W$  : additional measured adjustment variables, with  $F = F_0 \cup F_1$ .

$H \subseteq V \cup W$  : additional unobserved variables.

We will make use of the same extended SCM  $\tilde{M}$  with intervention variables  $I_b$  for  $b \in B$  and graph  $G_{\text{do}(I_B)}$  as for stating the do-calculus.

**Theorem 7.4.1** (General adjustment formula for simple SCMs). *Given a simple SCM  $M = (J, V, W, \mathcal{X}, P, f)$  with graph  $G = G(M)$ . Assume that all the following  $\sigma$ -separations hold in the graph  $G_{\text{do}(I_B, D)}$ :*

$$(F_0 \cup H) \underset{G_{\text{do}(I_B, D)}}{\perp}^{\sigma} I_B \mid (C \cup D), \quad (43)$$

$$A \underset{G_{\text{do}(I_B, D)}}{\perp}^{\sigma} (F_1 \cup I_B) \mid (B \cup F_0 \cup H \cup C \cup D), \quad (44)$$

$$H \underset{G_{\text{do}(I_B, D)}}{\perp}^{\sigma} B \mid (F \cup C \cup I_B \cup D). \quad (45)$$

Further assume that we have reference measures  $\mu_v$  on  $\mathcal{X}_v$ ,  $v \in V \cup H$ , such that:

$$\begin{aligned} \mu_{B \cup C \cup F \cup H} &\ll P(X_B, X_C, X_F, X_H \mid \text{do}(X_D)) \ll \mu_{B \cup C \cup F \cup H}, \\ \mu_{C \cup F \cup H} &\ll P(X_C, X_F, X_H \mid \text{do}(X_B, X_D)) \ll \mu_{C \cup F \cup H}. \end{aligned}$$

Then we have the adjustment formula:

$$P_M(X_A \mid X_C, \text{do}(X_B, X_D)) = P_M(X_A \mid X_B, X_C, X_F, \text{do}(X_D)) \circ P_M(X_F \mid X_C, \text{do}(X_D)) \quad \mu_{B \cup C}\text{-a.s.}$$

*Proof.* Analogous to that of Theorem 5.2.3, but now using the global Markov property for simple SCMs, Corollary 7.2.1, instead of the one for causal Bayesian networks, Theorem 4.2.1.  $\square$

The following is just the special case  $F_1 = H = \emptyset$ .

**Corollary 7.4.2** (Conditional interventional backdoor covariate adjustment formula for simple SCMs). *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G = G(M)$ . Assume that the conditional interventional backdoor criterion in the graph  $G_{\text{do}(I_B, D)}$  holds:*

1.  $F \underset{G_{\text{do}(I_B, D)}}{\perp}^{\sigma} I_B \mid (C \cup D)$ , and:
2.  $A \underset{G_{\text{do}(I_B, D)}}{\perp}^{\sigma} I_B \mid (B \cup F \cup C \cup D)$ .

Further assume the following absolute continuities:

$$\begin{aligned} \mu_{B \cup C \cup F} &\ll P(X_B, X_C, X_F \mid \text{do}(X_D)) \ll \mu_{B \cup C \cup F}, \\ \mu_{C \cup F} &\ll P(X_C, X_F \mid \text{do}(X_B, X_D)) \ll \mu_{C \cup F}. \end{aligned}$$

Then we have the adjustment formula:

$$P_M(X_A \mid X_C, \text{do}(X_B, X_D)) = P_M(X_A \mid X_B, X_C, X_F, \text{do}(X_D)) \circ P_M(X_F \mid X_C, \text{do}(X_D)) \quad \mu_{B \cup C}\text{-a.s.}$$

Without the conditioning set, i.e.  $C = \emptyset$ , and direct careful analysis we get a version with slightly weaker positivity assumptions:

**Theorem 7.4.3** (Interventional backdoor covariate adjustment formula). *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G = G(M)$ . Assume that the interventional backdoor criterion in the graph  $G_{\text{do}(I_B, D)}$  holds:*

1.  $F \perp_{G_{\text{do}(I_B, D)}}^{\sigma} I_B \mid D$ , and:
2.  $A \perp_{G_{\text{do}(I_B, D)}}^{\sigma} I_B \mid (B \cup F \cup D)$ .

Further assume the following absolute continuity:

$$P_M(X_F \mid \text{do}(X_D)) \otimes P_M(X_B \mid \text{do}(X_D)) \ll P_M(X_F, X_B \mid \text{do}(X_D)).$$

Then we have the adjustment formulas:

$$\begin{aligned} P_M(X_A, X_F \mid \text{do}(X_B, X_D)) &= P_M(X_A \mid X_F, X_B, \text{do}(X_D)) \otimes P_M(X_F \mid \text{do}(X_D)) \quad P_M(X_B \mid \text{do}(X_D))\text{-a.s.}, \\ P_M(X_A \mid \text{do}(X_B, X_D)) &= P_M(X_A \mid X_F, X_B, \text{do}(X_D)) \circ P_M(X_F \mid \text{do}(X_D)) \quad P_M(X_B \mid \text{do}(X_D))\text{-a.s.} \end{aligned}$$

*Proof.* Analogous to that of Corollary 5.2.5, but now using the global Markov property for simple SCMs, Corollary 7.2.1, instead of the one for causal Bayesian networks, Theorem 4.2.1.  $\square$

We can now further specialize to the case with  $C = D = J = \emptyset$  and immediately get:

**Corollary 7.4.4** (Backdoor covariate adjustment for simple SCMs). *Given a simple SCM  $M = (J, V, W, \mathcal{X}, P, f)$  with graph  $G = G(M)$ . Assume that the backdoor criterion holds:*

1.  $F \perp_{G_{\text{do}(I_B)}}^{\sigma} I_B$ , and:
2.  $A \perp_{G_{\text{do}(I_B)}}^{\sigma} I_B \mid (B \cup F)$ .

Further assume the following absolute continuity:

$$P_M(X_F) \otimes P_M(X_B) \ll P_M(X_F, X_B).$$

Then we have the adjustment formulae:

$$\begin{aligned} P_M(X_A, X_F \mid \text{do}(X_B)) &= P_M(X_A \mid X_F, X_B) \otimes P_M(X_F) && P_M(X_B)\text{-a.s.}, \\ P_M(X_A \mid \text{do}(X_B)) &= P_M(X_A \mid X_F, X_B) \circ P_M(X_F) && P_M(X_B)\text{-a.s.} \end{aligned}$$

The literature often fails to mention the strict positivity assumptions, even though without sufficient positivity, the various backdoor criteria may not hold. A simple example of how the adjustment formula may fail if the strict positivity assumptions are not met is provided in Example 5.3.29.

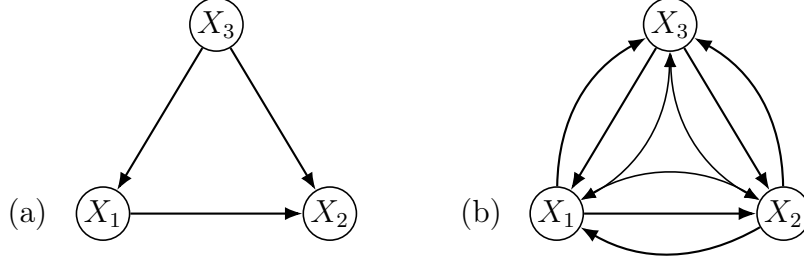


Figure 20: Two graphs of simple SCMs. In (a), the Markov kernel  $P_M(X_2 | \text{do}(X_1), X_3)$  is identifiable (under positivity assumptions) from  $P_M(X_1, X_2, X_3)$ . In (b), it is not identifiable, but we can still bound it using the natural bounds if both  $X_1$  and  $X_3$  are discrete.

## 7.5. Bounds on causal effects

If causal effects cannot be identified from the observable distribution, we may still be able to derive informative bounds (see also Figure 20). We first prove the *consistency* property of solution functions of simple SCMs.

**Proposition 7.5.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  be the solution function of  $M$ . Let  $V = A \dot{\cup} B$  be a partition of the endogenous variables of  $M$ , and let  $\tilde{g} : \mathcal{X}_A \times \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_B$  be the solution function of  $M_{\text{do}(X_A)}$ . Then*

$$g_B(x_J, x_W) = \tilde{g}(g_A(x_J, x_W), x_J, x_W)$$

for all  $x \in \mathcal{X}$ .

*Proof.* For all  $x \in \mathcal{X}$ :

$$\begin{cases} x_A = f_A(x), \\ x_B = f_B(x) \end{cases} \iff \begin{cases} x_A = g_A(x_J, x_W), \\ x_B = g_B(x_J, x_W). \end{cases}$$

Also, for all  $x \in \mathcal{X}$ :

$$x_B = f_B(x) \iff x_B = \tilde{g}_B(x_A, x_J, x_W).$$

Hence, for all  $x \in \mathcal{X}$ :

$$\begin{cases} x_A = f_A(x), \\ x_B = f_B(x) \end{cases} \implies \begin{cases} x_A = g_A(x_J, x_W), \\ x_B = \tilde{g}_B(x_A, x_J, x_W) \end{cases} \implies \begin{cases} x_A = g_A(x_J, x_W), \\ x_B = \tilde{g}_B(g_A(x_J, x_W), x_J, x_W). \end{cases}$$

The uniqueness of the solution function  $g$  of  $M$  now implies the consistency statement.  $\square$

[MN98] proved the following ‘natural’ bounds. We point out here that they also hold for simple SCMs.

**Theorem 7.5.2** (Natural bounds on causal effect). *Let  $M$  be a simple SCM with endogenous variables  $V \supseteq \{1, 2, 3\}$  and no exogenous input variables. Assume that  $\mathcal{X}_1$  is discrete. Then*

$$\begin{aligned} P_M(X_1 = a, X_2 \in B | X_3 \in C) &\leq P_M(X_2 \in B | X_3 \in C, \text{do}(X_1 = a)) \\ &\leq P_M(X_1 = a, X_2 \in B | X_3 \in C) + P_M(X_1 \neq a | X_3 \in C). \end{aligned} \quad (46)$$

for any  $a \in \mathcal{X}_1$ , measurable  $B \subseteq \mathcal{X}_2$ , and measurable  $C \subseteq \mathcal{X}_3$  with  $P_M(X_3 \in C) > 0$ .

*Proof.* We use the consistency  $X_2 = X_2^{\text{do}(X_1=X_1)}$  and elementary probability theory:

$$\begin{aligned} P_M(X_1 = a, X_2 \in B | X_3 \in C) &= P_M(X_1 = a, X_2^{\text{do}(X_1=a)} \in B | X_3 \in C) \\ &\leq P_M(X_2^{\text{do}(X_1=a)} \in B | X_3 \in C) \\ &= P_M(X_1 = a, X_2^{\text{do}(X_1=a)} \in B | X_3 \in C) \\ &\quad + P_M(X_1 \neq a, X_2^{\text{do}(X_1=a)} \in B | X_3 \in C) \\ &\leq P_M(X_1 = a, X_2 \in B | X_3 \in C) + P_M(X_1 \neq a | X_3 \in C) \end{aligned}$$

The statement follows since

$$P_M(X_2 \in B | X_3 \in C, \text{do}(X_1 = a)) = P_M(X_2^{\text{do}(X_1=a)} \in B | X_3 \in C).$$

□

This so-called “natural” bound can be shown to be tight. Remarkably, we do not need to make any assumptions regarding the causal relations between the three endogenous variables. This allows us to bound the causal effect of  $X_1$  on  $X_2$  in the presence of unobserved confounding, selection bias, and even cycles. Unfortunately, it can be shown that there exists no analogous bound in case  $X_1$  is real-valued.

If one has a priori knowledge about the range of  $X_2$  (in case it is real-valued), one can also derive a bound on the expected result of an intervention [MP13].

**Corollary 7.5.3.** *In the situation of Theorem 7.5.2, suppose that  $\mathcal{X}_2 = [\alpha, \beta]$  and  $0 < P_M(X_1 = a | X_3 \in C) < 1$ . Then*

$$\begin{aligned} &P_M(X_1 = a | X_3 \in C) \mathbb{E}_M(X_2 | X_1 = a, X_3 \in C) + \alpha P_M(X_1 \neq a | X_3 \in C) \\ &\leq \mathbb{E}_M(X_2 | X_3 \in C, \text{do}(X_1 = a)) \\ &\leq P_M(X_1 = a | X_3 \in C) \mathbb{E}_M(X_2 | X_1 = a, X_3 \in C) + \beta P_M(X_1 \neq a | X_3 \in C). \end{aligned} \quad (47)$$

*Proof.* Using consistency, we get:

$$\begin{aligned} &P_M(X_2^{\text{do}(X_1=a)} \in B | X_3 \in C) \\ &= P_M(X_1 = a | X_3 \in C) P_M(X_2^{\text{do}(X_1=a)} \in B | X_1 = a, X_3 \in C) \\ &\quad + P_M(X_1 \neq a | X_3 \in C) P_M(X_2^{\text{do}(X_1=a)} \in B | X_1 \neq a, X_3 \in C) \\ &= P_M(X_1 = a | X_3 \in C) P_M(X_2 \in B | X_1 = a, X_3 \in C) \\ &\quad + P_M(X_1 \neq a | X_3 \in C) P_M(X_2^{\text{do}(X_1=a)} \in B | X_1 \neq a, X_3 \in C). \end{aligned}$$

Integrating over  $X_2$ , the assumption  $\mathcal{X}_2 = [\alpha, \beta]$ , interval arithmetic, and affinity of expected values, gives:

$$\begin{aligned}
& \mathbb{E}_M(X_2^{\text{do}(X_1=a)} | X_3 \in C) \\
&= P_M(X_1 = a | X_3 \in C) \mathbb{E}_M(X_2 | X_1 = a, X_3 \in C) \\
&+ P_M(X_1 \neq a | X_3 \in C) \mathbb{E}_M(X_2^{\text{do}(X_1=a)} | X_1 \neq a, X_3 \in C) \\
&\in P_M(X_1 = a | X_3 \in C) \mathbb{E}_M(X_2 | X_1 = a, X_3 \in C) + P_M(X_1 \neq a | X_3 \in C) [\alpha, \beta].
\end{aligned}$$

□



## 8. Counterfactuals

Counterfactuals are questions of the kind “Fred obtained a cum laude for his PhD; would he have obtained it also if he were female?”. Counterfactuals consider a hypothetical situation that is “contrary to the fact”, that is, which differs from what was actually observed. One of the big practical obstacles of dealing with counterfactual probabilities is that they are typically not identifiable from experimental data, and at best only bounds on such quantities can be obtained. For systems with (almost) deterministic causal relations, these bounds may become quite informative, but tend to become more loose as the stochasticity in the system increases. While it would thus perhaps be easiest to avoid counterfactuals altogether, they do appear naturally in law and engineering. Humans also have a tendency to communicate using counterfactuals, and the grammar of many languages distinguishes counterfactual statements. This might be an inductive bias towards dealing with (almost) deterministic systems.

In this chapter, we will focus on counterfactuals that are hypothetical statements (or questions) regarding the effects of some action that is contrary-to-fact, closely following Pearl’s approach to counterfactuals. One can consider other types of counterfactuals as well, for example “backtracking” counterfactuals. Dealing with counterfactuals appears to be one of the least well-defined, but perhaps also most intriguing, aspects of causality.

### 8.1. Modeling counterfactuals via twinning

For example, suppose you are healthy but drank too much beer last night and now suffer from a hangover. A counterfactual statement is then: “If I had not drunk so much beer yesterday, I would feel much better now.” This statement invites one to imagine an alternative world in which everything is the same as in the actual world, with the sole difference that you did not drink beer last night. We can then use our causal model of the world to predict the consequences of this action (e.g., since you were in a healthy state and did not drink so much beer, you most likely will feel well in this alternative world).<sup>37</sup> This example already shows the ambiguity typically encountered in counterfactuals: if for you, not drinking beer means that you drink wine instead, then you may actually feel worse than if you had drunk beer.

Indeed, the truth value of such statements is often hard to determine in case the “world” is partially latent or not fully understood. When debugging a computer program, one makes heavy use of counterfactuals: “if I had put a minus sign there, then the output of my program would have been correct”. In case the full source code is available, it is in principle straightforward to work out whether such a statement is correct or not, but it becomes more difficult if the full source code is not available. It becomes even more challenging when the output of the computer program is stochastic or it

---

<sup>37</sup>The word “counterfactual” is also commonly used in the causal inference literature in a weaker sense, describing the potential outcome of a hypothetical situation that is not necessarily “contrary to the fact”. For example, some would refer to “If I drink too much beer tonight, I will have a hangover tomorrow” as a counterfactual statement as well. It is important to be aware of this to avoid possible confusion. We will only use the word “counterfactual” in its strong (contrary-to-fact) sense.

is impossible to reproduce its complete input. Generally, in situations where the full causal mechanism is unknown, or exogenous randomness is latent, counterfactuals may not be well-defined quantities. Nevertheless, counterfactual thinking is very common in humans, and toddlers already bombard their parents with counterfactual questions, presumably as a means for them to build internal causal models of the world.

The mathematical formalization of counterfactuals proposed by Pearl provides some clarification, but it also points out their inherent complexity and strong dependence on the chosen model.<sup>38</sup> It mimics the reasoning steps we mentally perform when thinking about counterfactuals by constructing a “(f)actual” world and a parallel “counterfactual” world, which is minimally different in some aspect. The crucial (and often untestable) assumption is that the exogenous random variables have the same *values* in both worlds. A good analogy here is that of two identical twins that share the same latent genetics. Before we give the definition, we will introduce some bookkeeping notation.

**Notation 8.1.1.** *Given an index set  $Z$  we define a primed copy  $Z' := \{z' : z \in Z\}$ , where each  $z'$  is a “primed” copy of  $z$  (distinguishable from  $z$  itself because of the attached prime symbol). We will also write  $(z')^\circ = z$  for  $z \in Z$ , where the superscript  $\circ$  removes the prime, i.e., it maps back to the original of the primed index.*

The following operation on SCMs (also known as the “twin-network approach” of [BP94b]) provides one possible way of modeling counterfactuals.<sup>39</sup>

**Definition 8.1.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. We define the twinned SCM of  $M$  as the SCM  $M^{\text{twin}} = (J^{\text{twin}} = J \dot{\cup} J', V^{\text{twin}} = V \dot{\cup} V', W, \mathcal{X}^{\text{twin}}, P, f^{\text{twin}})$  with  $J' = \{j' : j \in J\}$  a copy of  $J$  and  $V' = \{v' : v \in V\}$  a copy of  $V$ , the twinned domain given by*

$$\mathcal{X}^{\text{twin}} = \mathcal{X}_J \times \mathcal{X}_{J'} \times \mathcal{X}_V \times \mathcal{X}_{V'} \times \mathcal{X}_W$$

where  $\mathcal{X}_{j'} = \mathcal{X}_j$  for all  $j \in J$  and  $\mathcal{X}_{v'} = \mathcal{X}_v$  for all  $v \in V$ , and the twinned causal mechanism components given by

$$f_u^{\text{twin}}((x_J, x_{J'}), (x_V, x_{V'}), x_W) = \begin{cases} f_u(x_J, x_V, x_W) & u \in V, \\ f_{u^\circ}(x_{J'}, x_{V'}, x_W) & u \in V'. \end{cases}$$

The twinning operation is used to create copies of variables (so that in addition to the one in the factual world, we have its copy in the counterfactual world) that can have different values to describe contrary-to-fact situations. A specific choice in modeling counterfactuals in this way is the assumption that *all exogenous random variables* have

<sup>38</sup>Consider this a warning before attempting to predict counterfactual statements in a data-driven way, for example, using a neural network.

<sup>39</sup>The twinning operation can be applied to any SCM, but not to any L-CBN, as L-CBNs typically do not explicitly model latent random variables. Only for the subclass of L-CBNs that are in SCM form (see Definition 4.4.3), i.e., for which every node with at least one parent comes with a deterministic Markov kernel, could we define a twinning operation that is analogous to the one we define for SCMs.

the same value in the actual and in the counterfactual world. This is a very strong (and typically untestable) assumption.

The English language has a special grammatical construct to express counterfactuals: “If I had studied better, I would have passed the exam,” instead of “If I study better, I will pass the exam.” For the first statement, we first twin the SCM and then intervene on it, for the second, we just intervene on the SCM and there is no need for twinning.<sup>40</sup>

The twinning operation is well-defined in the following sense.

**Proposition 8.1.3.** *Twinning preserves equivalence:  $M \equiv \tilde{M} \implies M^{\text{twin}} \equiv \tilde{M}^{\text{twin}}$ .*

Furthermore, the twinning operation is compatible with marginalization.

**Proposition 8.1.4.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. For  $L \subseteq V$  such that  $M_{[L]}$  is uniquely solvable,*

$$(M_{\setminus L})^{\text{twin}} = (M^{\text{twin}})_{\setminus (L \cup L')}.$$

*Proof.* This follows by writing out the definitions and checking commutativity of the operations performed on the various components of the SCM tuple one-by-one.  $\square$

Hard interventions are compatible with the twinning operation, in the following sense:

**Proposition 8.1.5.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM.*

- For  $T \subseteq J \cup V$ ,  $x_T \in \mathcal{X}_T$ :

$$(M^{\text{twin}})_{\text{do}(X_T=x_T, X_{T'}=x_T)} = (M_{\text{do}(X_T=x_T)})^{\text{twin}}.$$

- For  $T \subseteq W$ ,  $Q_T \in \prod_{t \in T} \mathcal{P}(\mathcal{X}_t)$ :

$$(M^{\text{twin}})_{\text{do}(X_T \sim Q_T)} = (M_{\text{do}(X_T \sim Q_T)})^{\text{twin}}.$$

- For  $T \subseteq J \cup V$ :

$$(M^{\text{twin}})_{\text{do}(T, T')} = (M_{\text{do}(T)})^{\text{twin}}.$$

*Proof.* These properties all follow by writing out the definitions and checking commutativity of the operations performed on the various components of the SCM tuple one-by-one.  $\square$

Another important property of the twinning operation is that it preserves unique solvability and simplicity.

---

<sup>40</sup>Note that when considering two potential outcomes  $X^{\text{do}(x_J)}, X^{\text{do}(x'_J)}$  of SCM  $M$  for different inputs  $x_J, x'_J$  we do *not* necessarily assume that  $X_W^{\text{do}(x_J)} = X_W^{\text{do}(x'_J)}$ ; we only assume that they have the same distribution, that is,  $X_W^{\text{do}(x_J)} \sim X_W^{\text{do}(x'_J)}$ . This choice has been made to avoid introducing implicitly defined “cross-world” assumptions. In our opinion, it is better to explicitly introduce such counterfactuals with the twinning construction, because this enforces one to think about which variables are shared across potential worlds and which are copied (resampled or reevaluated).

**Lemma 8.1.6.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. Let  $T_1 \subseteq J \cup V$ ,  $T_2 \subseteq J' \cup V'$  and  $T_3 \subseteq W$ . If  $g_{\text{do}(T_1 \cup T_3)} : \mathcal{X}_{J \cup T_1 \cup T_3} \times \mathcal{X}_{W \setminus T_3} \rightarrow \mathcal{X}_{V \setminus T_1}$  is a solution function of  $M_{\text{do}(T_1 \cup T_3)}$  and  $g_{\text{do}(T_2^\circ \cup T_3)} : \mathcal{X}_{J \cup T_2^\circ \cup T_3} \times \mathcal{X}_{W \setminus T_3} \rightarrow \mathcal{X}_{V' \setminus T_2}$  is a solution function of  $M_{\text{do}(T_2^\circ \cup T_3)}$  then*

$$g^{\text{twin}} : \mathcal{X}_{(J \cup T_1 \cup T_3) \cup (J' \cup T_2)} \times \mathcal{X}_{W \setminus T_3} \rightarrow \mathcal{X}_{V \setminus T_1} \times \mathcal{X}_{V' \setminus T_2}$$

$$: (x_{(J \cup T_1 \cup T_3) \cup (J' \cup T_2)}, x_{W \setminus T_3}) \mapsto (g_{\text{do}(T_1 \cup T_3)}(x_{J \cup T_1 \cup T_3}, x_{W \setminus T_3}), g_{\text{do}(T_2^\circ \cup T_3)}(x_{J' \cup T_2 \cup T_3}, x_{W \setminus T_3}))$$

*is a solution function of  $(M^{\text{twin}})_{\text{do}(T_1 \cup T_2 \cup T_3)}$ . In case  $g_{\text{do}(T_1 \cup T_3)}$  and  $g_{\text{do}(T_2^\circ \cup T_3)}$  are unique,  $g^{\text{twin}}$  is also the unique solution function of  $(M^{\text{twin}})_{\text{do}(T_1 \cup T_2 \cup T_3)}$ .*

*Proof.* Let  $f^{\text{twin}}$  denote the causal mechanism of  $M^{\text{twin}}$ . For all  $x \in \mathcal{X}_J \times \mathcal{X}_{J'} \times \mathcal{X}_V \times \mathcal{X}_{V'} \times \mathcal{X}_W$ ,

$$x_{(V \cup V') \setminus (T_1 \cup T_2)} = f_{(V \cup V') \setminus (T_1 \cup T_2)}^{\text{twin}}(x)$$

$$\iff \begin{cases} x_{V \setminus T_1} &= f_{V \setminus T_1}(x_J, x_V, x_W) \\ x_{V' \setminus T_2} &= f_{(V' \setminus T_2)^\circ}(x_{J'}, x_{V'}, x_W) \end{cases}$$

$$\iff \begin{cases} x_{V \setminus T_1} &= g_{\text{do}(T_1 \cup T_3)}(x_{J \cup T_1 \cup T_3}, x_{W \setminus T_3}) \\ x_{V' \setminus T_2} &= g_{\text{do}(T_2^\circ \cup T_3)}(x_{J' \cup T_2 \cup T_3}, x_{W \setminus T_3}) \end{cases}$$

$$\iff x_{(V \cup V') \setminus (T_1 \cup T_2)} = g^{\text{twin}}(x_{J \cup T_1 \cup T_3}, x_{J' \cup T_2}, x_{W \setminus T_3})$$

In case  $g_{\text{do}(T_1 \cup T_3)}$  and  $g_{\text{do}(T_2^\circ \cup T_3)}$  are unique, the “ $\iff$ ” becomes an “ $\iff$ ”.  $\square$

**Proposition 8.1.7.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. If  $M$  is uniquely solvable, then  $M^{\text{twin}}$  is uniquely solvable. Furthermore, if  $M$  is simple, then  $M^{\text{twin}}$  is simple.*

*Proof.* The first statement follows directly from Lemma 8.1.6 by taking  $T_1 = T_2 = \emptyset$ . The second statement follows from Lemma 8.1.6 in combination with Theorem 6.5.2.  $\square$

We can also define a twinning operation on graphs. We will only define this for graphs without bidirected edges.

**Definition 8.1.8.** *Let  $G = (J, V \cup W, E)$  be a CDG such that  $\text{Pa}^G(W) = \emptyset$ . Write  $J' := \{j' : j \in J\}$  and  $V' := \{v' : v \in V\}$  for copies of  $J$  and  $V$ , respectively. The twinned graph  $G^{\text{twin}(J, V)}$  is defined as the CDG  $(J \dot{\cup} J', V \dot{\cup} V' \cup W, E^{\text{twin}})$  with directed edges*

$$E^{\text{twin}} := E \cup \{w \rightarrow v' : w \in W, v \in V, w \rightarrow v \in E\} \cup \{i' \rightarrow v' : i \in J \cup V, v \in V, i \rightarrow v \in E\}.$$

*In words, we copy the nodes  $J \cup V$  (but not the nodes  $W$ ) and copy the edges accordingly.*

The graphical and the SCM twinning operations are compatible:

**Proposition 8.1.9.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. Then*

$$G(M)^{\text{twin}(J \cup V)} = G(M^{\text{twin}}).$$

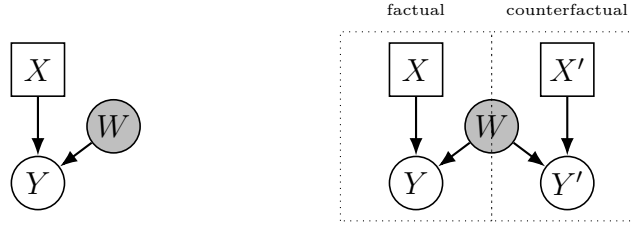


Figure 21: Left: Graph  $G(M)$  of SCM  $M$ . Right: Graph  $G(M^{\text{twin}}) = G(M)^{\text{twin}(X,Y)}$  of the twinned SCM  $M^{\text{twin}}$ . The exogenous random variable  $W$  is shared between the factual and counterfactual “world”, whereas the exogenous input variable  $X$  and the endogenous (output) variable  $Y$  may differ between the two worlds (with  $X'$  and  $Y'$  denoting the corresponding variables in the counterfactual “world”).

*Proof.* Follows by writing out the definitions. □

The simplest non-trivial example of a twin SCM is the following.

**Example 8.1.10.** Consider an SCM with exogenous input variable  $X$ , endogenous variable  $Y$ , exogenous random variable  $W$ , and structural equation:

$$Y^{\text{do}(x)} = f(x, W)$$

Its graph is depicted in Figure 21. Twinning adds an exogenous input variable  $X'$ , an endogenous variable  $Y'$ , and structural equation

$$(Y')^{\text{do}(x')} = f(x', W)$$

but keeps the same endogenous random variable  $W$ . The graph of the twinned SCM is also shown in Figure 21.

For instance,  $X$  could be the number of glasses of beer you consumed yesterday evening,  $Y$  the severity of a headache the next morning, and  $W$  would represent all other possible causes of a headache (e.g., COVID-19, a concussion obtained in a rugby game, the number of glasses of wine you consumed yesterday evening, ...). When stating “If I had not drunk so much beer yesterday, I would feel much better now,” we imagine a counterfactual world in which the value of  $W$  is the same as in the (f)actual world, but  $X'$  (and therefore  $Y'$ ) may have different values in the counterfactual world than the corresponding values  $X$  (and  $Y$ ) in the factual world.

We are not claiming that *all* counterfactuals can be modeled naturally using a twinned SCM.

**Example 8.1.11.** The counterfactual: “If I had not bumped into my old friend, I would not have ended up in the pub yesterday evening.” In this case, it appears most natural to model “bumping into an old friend” as an exogenous random variable, and it is precisely this exogenous random variable that is assumed to be different in the counterfactual world.

Obviously, we could model this by an SCM that has two copies of the exogenous random variable that represents “bumping into an old friend”.

## 8.2. Counterfactual Equivalence

In Definition 6.6.6, we defined the notions of observable and interventional equivalence for simple SCMs. We can add a more fine-grained notion of equivalence by making use of the twinning operation, which we refer to as counterfactual equivalence.<sup>41</sup>

**Definition 8.2.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  and  $\tilde{M} = (\tilde{J}, \tilde{V}, \tilde{W}, \tilde{\mathcal{X}}, \tilde{P}, \tilde{f})$  be two simple SCMs and  $O \subseteq (V \cup W) \cap (\tilde{V} \cup \tilde{W})$  a subset. We say that  $M$  and  $\tilde{M}$  are counterfactually equivalent w.r.t.  $O$  if the twin SCMs  $M^{\text{twin}}$  and  $\tilde{M}^{\text{twin}}$  are interventionally equivalent w.r.t.  $O \cup O'$ , where  $O'$  is the copy of  $O \cap (V \cap \tilde{V})$  in  $V' \cap \tilde{V}'$ .*

More generally, one could define counterfactual equivalence not only with respect to an observed set of variables, but also with respect to a given set of interventions.

We get the following important corollary of Theorem 6.7.8.

**Corollary 8.2.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM,  $L \subseteq V$ , and  $M_{\setminus L}$  its marginalization over  $L$ . Then  $M$  and  $M_{\setminus L}$  are observably, interventionally and counterfactually equivalent w.r.t.  $(V \cup W) \setminus L$ .*

*Proof.* The observable and interventional equivalence is the claim of Theorem 6.7.8. Write  $K = V \setminus L$ . By Proposition 8.1.4,  $(M_{\setminus L})^{\text{twin}} = (M^{\text{twin}})_{\setminus (L \cup L')}$ . Since  $M^{\text{twin}}$  and its marginalization  $(M^{\text{twin}})_{\setminus (L \cup L')}$  are interventionally equivalent w.r.t.  $K \cup K' \cup W$ ,  $M$  and  $M_{\setminus L}$  are counterfactually equivalent w.r.t.  $K \cup W$ .  $\square$

**Lemma 8.2.3.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Then  $M$  and  $M^{\text{twin}}$  are interventionally equivalent w.r.t.  $V \cup W$ .*

*Proof.* We have to show that for any  $T \subseteq V \cup W$ ,  $M_{\text{do}(T)}$  and  $(M^{\text{twin}})_{\text{do}(T)}$  are observably equivalent w.r.t.  $(V \cup W) \setminus T$ . Let  $T \subseteq V \cup W$ , and write  $T_1 = T \cap V$ ,  $T_2 = \emptyset$ ,  $T_3 = T \cap W$ . Then  $T = T_1 \cup T_3$ .

By Lemma 8.1.6, with  $g_{\text{do}(T)} : \mathcal{X}_{J \cup T} \times \mathcal{X}_{W \setminus T_3} \rightarrow \mathcal{X}_{V \setminus T_1}$  the solution function of  $M_{\text{do}(T)}$  and  $g_{\text{do}(T_3)} : \mathcal{X}_{J \cup T_3} \times \mathcal{X}_{W \setminus T_3} \rightarrow \mathcal{X}_V$  the solution function of  $M_{\text{do}(T_3)}$ ,

$$\begin{aligned} \tilde{g} : \mathcal{X}_{(J \cup T) \cup J'} \times \mathcal{X}_{W \setminus T_3} &\rightarrow \mathcal{X}_{V \setminus T_1} \times \mathcal{X}_{V'} \\ &: (x_{(J \cup T) \cup J'}, x_{W \setminus T_3}) \mapsto (g_{\text{do}(T)}(x_{J \cup T}, x_{W \setminus T_3}), g_{\text{do}(T_3)}(x_{J' \cup T_3}, x_{W \setminus T_3})) \end{aligned}$$

is the unique solution function of  $(M^{\text{twin}})_{\text{do}(T)}$ . Note that  $g_{\text{do}(T)} \circ \text{pr}_{\mathcal{X}_{J \cup T} \times \mathcal{X}_{W \setminus T_3}} = \tilde{g}_{V \setminus T}$ .

Then  $P_{M_{\text{do}(T)}}(X_{(V \cup W) \setminus T} \mid \text{do}(X_{J \cup T}))$  is the push-forward  $(g_{\text{do}(T)}, \text{id}_{\mathcal{X}_{W \setminus T}})_*(P_{W \setminus T})$  of the marginal exogenous distribution  $P_{W \setminus T}$  of  $M_{\text{do}(T)}$  (interpreted as a constant Markov kernel  $\mathcal{X}_{J \cup T} \dashrightarrow \mathcal{X}_{W \setminus T}$ ). We obtain the marginal Markov kernel  $P_{(M^{\text{twin}})_{\text{do}(T)}}(X_{(V \cup W) \setminus T} \mid \text{do}(X_{J \cup T \cup J'}))$  as the push-forward  $(\tilde{g}_{V \setminus T}, \text{id}_{\mathcal{X}_{W \setminus T}})_*(P_{W \setminus T})$  of the marginal exogenous distribution  $P_{W \setminus T}$  of  $(M^{\text{twin}})_{\text{do}(T)}$ , now interpreted as a constant Markov kernel  $\mathcal{X}_{J \cup T} \times \mathcal{X}_{J'} \rightarrow \mathcal{X}_{W \setminus T}$ . Since  $g_{\text{do}(T)} \circ \text{pr}_{\mathcal{X}_{J \cup T} \times \mathcal{X}_{W \setminus T_3}} = \tilde{g}_{V \setminus T}$ , we obtain the desired conclusion.  $\square$

<sup>41</sup>The definition of counterfactual equivalence for (possibly non-simple) SCMs is provided in [BFPM21] for SCMs without exogenous input variables.

One can show the following properties of these equivalences:

**Proposition 8.2.4.** *For simple SCMs  $M, \tilde{M}$  and a subset  $O \subseteq (V \cap \tilde{V}) \cup (W \cap \tilde{W})$ :*

1. *If  $M \equiv \tilde{M}$  then  $M$  and  $\tilde{M}$  are counterfactually equivalent w.r.t.  $O$ .*
2. *If  $M$  and  $\tilde{M}$  are counterfactually equivalent w.r.t.  $O$  then  $M$  and  $\tilde{M}$  are interventionally equivalent w.r.t.  $O$ .*

*Proof.* 1. Suppose  $M \equiv \tilde{M}$ . Then  $M^{\text{twin}} \equiv \tilde{M}^{\text{twin}}$ , and the claim directly follows from Proposition 6.6.7.

2. Let  $O'$  be the copy of  $O \cap (V \cap \tilde{V})$  in  $V' \cap \tilde{V}'$ . Suppose  $M$  and  $\tilde{M}$  are counterfactually equivalent w.r.t.  $O$ . Then  $M^{\text{twin}}$  and  $\tilde{M}^{\text{twin}}$  are interventionally equivalent w.r.t.  $O \cup O'$ . For every  $T \subseteq O \cup O'$ ,  $(M^{\text{twin}})_{\text{do}(T)}$  and  $(\tilde{M}^{\text{twin}})_{\text{do}(T)}$  are observably equivalent w.r.t.  $(O \cup O') \setminus T$ .

We have to show that  $M$  and  $\tilde{M}$  are interventionally equivalent w.r.t.  $O$ . That is, we have to show that for every  $S \subseteq O$ ,  $M_{\text{do}(S)}$  and  $\tilde{M}_{\text{do}(S)}$  are observably equivalent w.r.t.  $O \setminus S$ .

Now take  $S \subseteq O$  and then let  $T_1 = S \cap V \cap \tilde{V}$ ,  $T_2 = T_1'$  and  $T_3 = S \cap W \cap \tilde{W}$ ; then  $S = T_1 \cup T_3$ . Then by assumption,

$$((M_{\text{do}(T_1)})^{\text{twin}})_{\text{do}(T_3)} = ((M^{\text{twin}})_{\text{do}(T_1, T_1')})_{\text{do}(T_3)} = (M^{\text{twin}})_{\text{do}(T_1, T_2, T_3)}$$

and

$$((\tilde{M}_{\text{do}(T_1)})^{\text{twin}})_{\text{do}(T_3)} = ((\tilde{M}^{\text{twin}})_{\text{do}(T_1, T_1')})_{\text{do}(T_3)} = (\tilde{M}^{\text{twin}})_{\text{do}(T_1, T_2, T_3)}$$

are observably equivalent w.r.t.  $(O \cup O') \setminus T = (O \setminus (T_1 \cup T_3)) \cup (O' \setminus T_2)$ , and hence w.r.t.  $O \setminus S = O \setminus (T_1 \cup T_3)$ . By Lemma 8.2.3,  $(M_{\text{do}(T_1)})^{\text{twin}}$  and  $M_{\text{do}(T_1)}$  are interventionally equivalent w.r.t.  $V \setminus T_1 \cup W$ . Hence  $((M_{\text{do}(T_1)})^{\text{twin}})_{\text{do}(T_3)}$  and  $(M_{\text{do}(T_1)})_{\text{do}(T_3)} = M_{\text{do}(S)}$  are observably equivalent w.r.t.  $V \setminus T_1 \cup W \setminus T_3$ , and hence also w.r.t.  $O \setminus S$ . Similarly,  $((\tilde{M}_{\text{do}(T_1)})^{\text{twin}})_{\text{do}(T_3)}$  and  $(\tilde{M}_{\text{do}(T_1)})_{\text{do}(T_3)} = \tilde{M}_{\text{do}(S)}$  are observably equivalent w.r.t.  $O \setminus S$ . Hence, by transitivity,  $M_{\text{do}(S)}$  and  $\tilde{M}_{\text{do}(S)}$  are observably equivalent w.r.t.  $O \setminus S$ . □

However, the reverse implications do not hold in general. Together with Proposition 6.6.7, Proposition 8.2.4 expresses that causal modeling is more refined than probabilistic modeling, and counterfactual modeling is more refined than interventional modeling. This formalizes what Pearl refers to as the “causal hierarchy” or “ladder of causation”.

In general, interventional equivalence does not imply counterfactual equivalence. Even interventionally equivalent SCMs with the same causal mechanism (that differ only in terms of their exogenous distributions) may not be counterfactually equivalent. For example, the SCMs  $M_\rho$  and  $M_{\rho'}$  with  $\rho \neq \rho'$  in the following example are interventionally equivalent, but not counterfactually equivalent.

**Example 8.2.5** (Interventional equivalence does not imply counterfactual equivalence [Daw02]). For parameter  $\rho \in [0, 1]$ , consider the SCM  $M_\rho$  with binary exogenous input variable  $X \in \{0, 1\}$ , endogenous variable  $Y \in \mathbb{R}$ , a single latent exogenous random variable  $W = (W_1, W_2) \in \mathbb{R}^2$  with exogenous distribution

$$\begin{pmatrix} W_1 \\ W_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$$

and structural equation

$$Y = W_1(1 - X) + W_2X.$$

In a medical setting, this SCM could be used to model whether a patient was treated or not ( $x = 1$  vs.  $x = 0$ ) and the corresponding potential outcome ( $Y^{\text{do}(x=1)}$  vs.  $Y^{\text{do}(x=0)}$ ).

Suppose that in the actual world we did not assign treatment to a patient ( $x = 0$ ) and the outcome was  $Y^{\text{do}(x=0)} = y \in \mathbb{R}$ . Consider the counterfactual query “What would the outcome have been, had we assigned treatment to this patient?”. We can answer this question by introducing a parallel counterfactual world in which the exogenous random variables for each patient have the same values as in the actual world, but treatment and outcome may differ. For this, consider the twin SCM  $M_\rho^{\text{twin}}$ . The counterfactual query then asks for

$$P_{M_\rho^{\text{twin}}}((Y')^{\text{do}(x'=1)} \mid Y^{\text{do}(x=0)} = y),$$

where  $Y^{\text{do}(x=0)}$  is the factual outcome, and  $(Y')^{\text{do}(x'=1)}$  is the counterfactual outcome (which are both marginal potential outcomes of the twinned SCM). One can calculate that

$$P_{M_\rho^{\text{twin}}}((Y')^{\text{do}(x'=1)}, Y^{\text{do}(x=0)}) = \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right)$$

and hence  $P_{M_\rho^{\text{twin}}}((Y')^{\text{do}(x'=1)} \mid Y^{\text{do}(x=0)} = y) = \mathcal{N}(\rho y, 1 - \rho^2)$  (by the general formula for conditioning a multivariate Gaussian distribution). Note that the answer to the counterfactual query depends on a quantity  $\rho$  that we cannot identify from the Markov kernel  $P_{M_\rho}(Y \mid \text{do}(X))$ , as it is independent of  $\rho$ . Therefore, even unlimited data from a randomized controlled trial would not suffice to determine the value of this particular counterfactual query. Indeed, SCMs  $M_\rho$  and  $M_{\rho'}$  with  $\rho \neq \rho'$  are interventionally equivalent, but not counterfactually equivalent.

The lesson of this example is that if one attempts to learn an SCM from data (even from randomized controlled trials with arbitrarily large sample size) it can happen that one still cannot identify the values of some counterfactual probabilities. In other words, data-driven estimation of counterfactual probabilities can be an ill-posed problem. Nevertheless, counterfactual are central in court cases (e.g., to determine responsibility, “the physician treated the patient with drug A and the patient died, would the patient still be alive if the physician had abstained from the treatment?”). The above example shows that one can be on very slippery terrain when it comes to answering such questions.



### 8.3. Exogenous reparameterizations

When exogenous random variables are latent, certain reparameterizations of those may preserve part of the causal semantics of the observed variables.

**Definition 8.3.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM. Let  $\tilde{W}$  be a finite index set disjoint from  $J \cup V \cup W$ ,  $\mathcal{X}_{\tilde{W}} = \prod_{\tilde{w} \in \tilde{W}} \mathcal{X}_{\tilde{w}}$  the product of standard measurable spaces  $\mathcal{X}_{\tilde{w}}$ . Let  $\tilde{f} : \mathcal{X}_J \times \mathcal{X}_V \times \mathcal{X}_{\tilde{W}} \rightarrow \mathcal{X}_V$  and  $\Phi : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_{\tilde{W}}$  be measurable mappings such that:*

i) *for all  $x_J \in \mathcal{X}_J, x_V \in \mathcal{X}_V, x_W \in \mathcal{X}_W$ ,*

$$f(x_J, x_V, x_W) = \tilde{f}(x_J, x_V, \Phi(x_J, x_W)),$$

*and*

ii) *the push-forward*

$$\tilde{P} := (\text{id}_{\mathcal{X}_J}, \Phi)_* \left( \delta(X_J | X_J) \otimes \bigotimes_{w \in W} P(X_w) \right) = \delta(X_J | X_J) \otimes \bigotimes_{\tilde{w} \in \tilde{W}} \tilde{P}_{\tilde{w}}$$

*factorizes with  $\tilde{P}_{\tilde{w}} \in \mathcal{P}(\mathcal{X}_{\tilde{w}})$  for  $\tilde{w} \in \tilde{W}$ .*

*Then we call the SCM*

$$M_{\text{repar}(\tilde{f}, \Phi)} = (J, V, \tilde{W}, \mathcal{X}_J \times \mathcal{X}_V \times \mathcal{X}_{\tilde{W}}, \tilde{P}, \tilde{f})$$

*an exogenous reparameterization of  $M$ .*

Exogenous reparameterizations are compatible with certain hard interventions that do *not* target the exogenous random variables.

**Lemma 8.3.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM, and  $M_{\text{repar}(\tilde{f}, \Phi)} = (J, V, \tilde{W}, \tilde{\mathcal{X}}, \tilde{P}, \tilde{f})$  an exogenous reparameterization of  $M$ . For an intervention target  $T \subseteq V \cup J$ :*

$$(M_{\text{repar}(\tilde{f}, \Phi)})_{\text{do}(T)} = (M_{\text{do}(T)})_{\text{repar}(\tilde{f}_{\setminus T}, \Phi)}.$$

*Proof.* This follows by writing out the definitions. □

The twinning operation is only compatible with exogenous reparameterizations under restrictive assumptions.

**Lemma 8.3.3.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM, and  $M_{\text{repar}(\tilde{f}, \Phi)} = (J, V, \tilde{W}, \tilde{\mathcal{X}}, \tilde{P}, \tilde{f})$  an exogenous reparameterization of  $M$ . If  $\Phi$  is constant in  $X_J$ , that is,  $\Phi(x_J, x_W) = \Phi(\tilde{x}_J, x_W)$  for all  $x_W \in \mathcal{X}_W$  and all  $x_J, \tilde{x}_J \in \mathcal{X}_J$ , then:*

$$(M_{\text{repar}(\tilde{f}, \Phi)})^{\text{twin}} = (M^{\text{twin}})_{\text{repar}(\tilde{f}^{\text{twin}}, \Phi)}.$$

*Proof.* This follows by writing out the definitions. In particular, we check that

$$\begin{aligned} f^{\text{twin}}(x_J, x_{J'}, x_V, x_{V'}, x_W) &= (f(x_J, x_V, x_W), f(x_{J'}, x_{V'}, x_W)) \\ &= (\tilde{f}(x_J, x_V, \Phi(x_W)), \tilde{f}(x_{J'}, x_{V'}, \Phi(x_W))) \\ &= \tilde{f}^{\text{twin}}(x_J, x_{J'}, x_V, x_{V'}, \Phi(x_W)). \end{aligned}$$

For the last equality, we need that  $\Phi$  does not depend on  $X_J$ .  $\square$

**Theorem 8.3.4.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be an SCM, and  $M_{\text{repar}(\tilde{f}, \Phi)} = (J, V, \tilde{W}, \tilde{\mathcal{X}}, \tilde{P}, \tilde{f})$  an exogenous reparameterization of  $M$ . If both  $M$  and  $M_{\text{repar}(\tilde{f}, \Phi)}$  are simple,<sup>42</sup> then  $M$  is observably and interventionally equivalent to  $M_{\text{repar}(\tilde{f}, \Phi)}$  w.r.t.  $V$ . If, furthermore,  $\Phi$  does not depend on  $X_J$ , then  $M$  is even counterfactually equivalent to  $M_{\text{repar}(\tilde{f}, \Phi)}$  w.r.t.  $V$ .*

*Proof.* Let  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  be the solution function of  $M$ , and  $\tilde{g} : \mathcal{X}_J \times \mathcal{X}_{\tilde{W}} \rightarrow \mathcal{X}_V$  the solution function of  $\tilde{M} := M_{\text{repar}(\tilde{f}, \Phi)}$ .  $g$  satisfies

$$x_V = f(x_J, x_V, x_W) \iff x_V = g(x_J, x_W)$$

for all  $x_J \in \mathcal{X}_J, x_V \in \mathcal{X}_V$  and  $x_W \in \mathcal{X}_W$ , while  $\tilde{g}$  satisfies

$$x_V = \tilde{f}(x_J, x_V, x_{\tilde{W}}) \iff x_V = \tilde{g}(x_J, x_{\tilde{W}})$$

for all  $x_J \in \mathcal{X}_J, x_V \in \mathcal{X}_V$  and  $x_{\tilde{W}} \in \mathcal{X}_{\tilde{W}}$ . By assumption, for all  $x_J \in \mathcal{X}_J, x_V \in \mathcal{X}_V, x_W \in \mathcal{X}_W$ ,

$$f(x_J, x_V, x_W) = \tilde{f}(x_J, x_V, \Phi(x_J, x_W)).$$

Hence, the following equation holds for all  $x_J \in \mathcal{X}_J, x_W \in \mathcal{X}_W$ :

$$g(x_J, x_W) = \tilde{g}(x_J, \Phi(x_J, x_W)).$$

In other words,

$$g = \tilde{g} \circ (\text{id}_{\mathcal{X}_J}, \Phi).$$

We now show that the marginal Markov kernels  $P_M(X_V \mid \text{do}(X_J))$  and  $P_{M_{\text{repar}(\tilde{f}, \Phi)}}(X_V \mid \text{do}(X_J))$  are the same. By Theorem 6.5.2,

$$\begin{aligned} P_M(X_V \mid \text{do}(X_J)) &= (g)_* (\delta(X_J \mid X_J) \otimes P(X_W)) \\ &= (\tilde{g} \circ (\text{id}_{\mathcal{X}_J}, \Phi))_* (\delta(X_J \mid X_J) \otimes P(X_W)) \\ &= (\tilde{g})_* (\delta(X_J \mid X_J) \otimes \tilde{P}(X_{\tilde{W}})) \\ &= P_{M_{\text{repar}(\tilde{f}, \Phi)}}(X_V \mid \text{do}(X_J)). \end{aligned}$$

This shows the observable equivalence w.r.t.  $V$ .

<sup>42</sup>This assumption is only needed because we avoided to define the equivalence relations for arbitrary SCMs.

Let  $T \subseteq V$ . Then  $(M_{\text{repar}(\tilde{f}, \Phi)})_{\text{do}(T)} = (M_{\text{do}(T)})_{\text{repar}(\tilde{f}_{\setminus T}, \Phi)}$  by Lemma 8.3.2. The observable equivalence of  $M_{\text{do}(T)}$  and  $(M_{\text{do}(T)})_{\text{repar}(\tilde{f}_{\setminus T}, \Phi)}$  w.r.t.  $V \setminus T$  hence implies the observable equivalence of  $M_{\text{do}(T)}$  and  $(M_{\text{repar}(\tilde{f}, \Phi)})_{\text{do}(T)}$  w.r.t.  $V \setminus T$ . Since this holds for all  $T \subseteq V$ ,  $M$  and  $M_{\Phi}$  are interventionally equivalent w.r.t.  $V$ .

Assume now that  $\Phi$  does not depend on  $X_J$ . By Lemma 8.3.3,  $(M_{\text{repar}(\tilde{f}, \Phi)})^{\text{twin}} = (M^{\text{twin}})_{\text{repar}(\tilde{f}^{\text{twin}}, \Phi)}$ . Since  $M^{\text{twin}}$  and its exogenous reparameterization  $(M^{\text{twin}})_{\text{repar}(\tilde{f}^{\text{twin}}, \Phi)}$  are interventionally equivalent w.r.t.  $V \cup V'$ ,  $M$  and  $M_{\text{repar}(\tilde{f}, \Phi)}$  are counterfactually equivalent w.r.t.  $V$ .  $\square$

A special case of interest is obtained for ‘pointwise’ surjective mappings  $\Phi$ .

**Corollary 8.3.5.** *If  $x_W \mapsto \Phi(x_J, x_W)$  is surjective for all  $x_J \in \mathcal{X}_J$ , then unique solvability of  $M$  implies unique solvability of  $M_{\text{repar}(\tilde{f}, \Phi)}$ , and simplicity of  $M$  implies simplicity of  $M_{\text{repar}(\tilde{f}, \Phi)}$ .*

*Proof.* Let  $g : \mathcal{X}_J \times \mathcal{X}_W \rightarrow \mathcal{X}_V$  be the solution function of  $M$ . It satisfies

$$x_V = f(x_J, x_V, x_W) \iff x_V = g(x_J, x_W)$$

for all  $x_J \in \mathcal{X}_J, x_V \in \mathcal{X}_V$  and  $x_W \in \mathcal{X}_W$ . Let  $\Psi : \mathcal{X}_J \times \mathcal{X}_{\tilde{W}} \rightarrow \mathcal{X}_W$  be an arbitrary left-inverse of  $\Phi$  (i.e.,  $\Phi(x_J, \Psi(x_J, x_{\tilde{W}})) = x_{\tilde{W}}$  for all  $x_J \in \mathcal{X}_J, x_{\tilde{W}} \in \mathcal{X}_{\tilde{W}}$ , implying  $(\text{id}_{\mathcal{X}_J}, \Phi) \circ (\text{id}_{\mathcal{X}_J}, \Psi) = (\text{id}_{\mathcal{X}_J}, \text{id}_{\mathcal{X}_{\tilde{W}}})$ ) and define

$$\tilde{g}(x_J, x_{\tilde{W}}) := g(x_J, \Psi(x_J, x_{\tilde{W}}))$$

for all  $x_J \in \mathcal{X}_J, x_{\tilde{W}} \in \mathcal{X}_{\tilde{W}}$ . Then

$$\begin{aligned} x_V = \tilde{f}(x_J, x_V, x_{\tilde{W}}) &\iff x_V = \tilde{f}(x_J, x_V, \Phi(x_J, \Psi(x_J, x_{\tilde{W}}))) \\ &\iff x_V = f(x_J, x_V, \Psi(x_J, x_{\tilde{W}})) \\ &\iff x_V = g(x_J, \Psi(x_J, x_{\tilde{W}})) \\ &\iff x_V = \tilde{g}(x_J, x_{\tilde{W}}) \end{aligned}$$

for all  $x_J \in \mathcal{X}_J, x_V \in \mathcal{X}_V, x_{\tilde{W}} \in \mathcal{X}_{\tilde{W}}$ . So  $\tilde{g}$  is the unique solution function of  $M_{\text{repar}(\tilde{f}, \Phi)}$ . Hence, unique solvability of  $M$  implies unique solvability of  $M_{\text{repar}(\tilde{f}, \Phi)}$  by Theorem 6.5.2.

Now suppose that  $M$  is simple. Let  $T \subseteq V$ . By Lemma 8.3.2,

$$(M_{\text{repar}(\tilde{f}, \Phi)})_{\text{do}(T)} = (M_{\text{do}(T)})_{\text{repar}(\tilde{f}_{\setminus T}, \Phi)}.$$

The unique solvability of  $M_{\text{do}(T)}$  implies that of  $(M_{\text{do}(T)})_{\text{repar}(\tilde{f}_{\setminus T}, \Phi)}$ , and hence that of  $(M_{\text{repar}(\tilde{f}, \Phi)})_{\text{do}(T)}$ . Therefore (with Remark 6.5.6),  $M_{\text{repar}(\tilde{f}, \Phi)}$  is simple.  $\square$

The following example shows that an exogenous reparameterization need not be counterfactually equivalent w.r.t.  $V$  if  $\Phi$  depends on  $X_J$ .

**Example 8.3.6.** Consider the acyclic SCM  $M$  with exogenous input variable  $X_1$  with co-domain  $\{-1, 1\}$ , endogenous variables  $X_2, X_3$  with co-domains  $\{-1, 1\}$ ,  $\{-2, 0, 2\}$ , respectively, and causal mechanism

$$\begin{aligned} X_2 &= f_2(X_1, X_B) = X_1 X_B \\ X_3 &= f_3(X_2, X_B) = X_2 + X_B, \end{aligned}$$

with exogenous random variable  $X_B \sim \text{Uni}(\{-1, 1\})$ . Consider the exogenous reparameterization  $\tilde{M}$  with exogenous random variable  $X_{\tilde{B}}$  with co-domain  $\{-1, 1\}$ , mapping  $\Phi : \{-1, 1\}^2 \rightarrow \{-1, 1\} : (x_1, x_B) \mapsto x_1 x_B$ , and causal mechanism

$$\begin{aligned} X_2 &= \tilde{f}_2(X_{\tilde{B}}) = X_{\tilde{B}} \\ X_3 &= \tilde{f}_3(X_1, X_2, X_{\tilde{B}}) = X_2 + X_1 X_{\tilde{B}}. \end{aligned}$$

Its exogenous distribution is  $\Phi_*(\mathbb{P}(X_B)) = \text{Uni}(\{-1, 1\}) = \mathbb{P}(X_{\tilde{B}})$ . It is indeed an exogenous reparameterization:

$$\begin{aligned} \tilde{f}(x_1, x_2, \Phi(x_1, x_B)) &= (\Phi(x_1, x_B), x_2 + x_1 \Phi(x_1, x_B)) \\ &= (x_1 x_B, x_2 + x_1 x_1 x_B) \\ &= (x_1 x_B, x_2 + x_B) \\ &= f(x_1, x_2, x_B). \end{aligned}$$

Both  $M$  and  $\tilde{M}$  are acyclic, hence simple. By Theorem 8.3.4,  $\tilde{M}$  is observably and interventionally equivalent to  $M$  w.r.t.  $\{X_2, X_3\}$ . However,  $\tilde{M}$  is not counterfactually equivalent to  $M$  w.r.t.  $\{X_2, X_3\}$ , as one can check explicitly.

## 8.4. Parameterizing SCMs using response functions

The following technique of “response functions” [BP94a] provides an example of an exogenous reparameterization. It has been used—amongst others—to derive bounds on counterfactual probabilities and to obtain tests for valid instruments. Here we will explain the idea using an example rather than with a general formal treatment (in that way avoiding some heavy bookkeeping).

**Definition 8.4.1.** Let  $M = (J, V = \{v\}, W, \mathcal{X}_J \times \mathcal{X}_V \times \mathcal{X}_W, P, f)$  be a simple SCM with  $\mathcal{X}_J$  discrete,  $\mathcal{X}_v$  discrete, and  $\mathcal{X}_W$  an arbitrary standard measurable space.

Let  $\tilde{W} = \{\tilde{w}\}$  be a singleton and consider the space of (measurable) functions<sup>43</sup>

$$\mathcal{X}_{\tilde{W}} := (\mathcal{X}_v)^{\mathcal{X}_J} = \{\phi : \mathcal{X}_J \rightarrow \mathcal{X}_v\}.$$

The SCM  $M$  induces a function  $\Phi : \mathcal{X}_W \rightarrow \mathcal{X}_{\tilde{W}}$  that assigns to each exogenous random value  $x_W \in \mathcal{X}_W$  the corresponding response function in  $\mathcal{X}_{\tilde{W}}$ , that is,  $\Phi(x_W)$  is the function  $x_J \mapsto f_v(x_J, x_W)$ . The SCM  $M_{\text{repar}(\tilde{f}, \Phi)} = (J, V = \{v\}, \tilde{W}, \mathcal{X}_J \times \mathcal{X}_V \times \mathcal{X}_{\tilde{W}}, \tilde{P}, \tilde{f})$

<sup>43</sup> [BP94a] call these functions “response functions”, and a random variable taking values in  $\mathcal{X}_{\tilde{W}}$  a “response-function variable”.

with an exogenous distribution the push-forward  $\tilde{P} = (\Phi)_*(P)$  and causal mechanism  $\tilde{f}(x_J, x_{\tilde{W}}) = x_{\tilde{W}}(x_J)$  (which just evaluates the response function  $x_{\tilde{W}}$  in the input  $x_J$ ) is called a response variable parameterization of  $M$ .

We call it a parameterization because it preserves the important causal semantics:

**Proposition 8.4.2.** *The response variable parameterization of  $M$  in Definition 8.4.1 is an exogenous reparameterization of  $M$  and is counterfactually equivalent to  $M$  w.r.t.  $V$ .*

*Proof.* Note that:

$$f(x_J, x_w) = \Phi(x_w)(x_J) = \tilde{f}(x_J, \Phi(x_w))$$

for all  $x_J \in \mathcal{X}_J, x_w \in \mathcal{X}_w$ . Because  $M$  and  $\tilde{M}$  are both acyclic, they are both simple, and the claim now follows from Theorem 8.3.4, noting that  $\Phi$  does not depend on  $X_J$ .  $\square$

**Corollary 8.4.3.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with discrete exogenous input space  $\mathcal{X}_J$ . Let  $v \in V$  be an endogenous variable in  $M$  taking values in a discrete space  $\mathcal{X}_v$ . Then there exists an SCM that is counterfactually equivalent to  $M$  w.r.t.  $\{v\}$  that has just a single endogenous variable ( $V = \{v\}$  with space  $\mathcal{X}_v$ ) and exogenous random space  $(\mathcal{X}_v)^{\mathcal{X}_J}$ .*

*Proof.* First marginalize out all endogenous variables except  $v$ , and then take the response variable parameterization of this marginal SCM. It has the desired exogenous random space, and note that both operations preserve counterfactual equivalence w.r.t.  $\{v\}$ .  $\square$

## 8.5. Bounding counterfactual probabilities

We have seen that unless one is willing to make very strong modeling assumptions, obtaining counterfactual probabilities from data can be impossible. In certain cases, though, it is possible to derive *bounds* on counterfactual probabilities from (intervened) Markov kernels [BP94a]. One way to derive such bounds on counterfactuals exploits the response function parameterization.

As a motivation, consider a medical setting in which a patient may either be treated (or not) and a week later the patient is cured (or not). Suppose a patient participating in a randomized controlled trial was assigned to the control group and hence not treated ( $\text{do}(x = 0)$ ), and it turned out one week later that this patient was not cured ( $Y^{\text{do}(x=0)} = 0$ ). The patient now wonders “would I have been cured, had I been assigned to the treatment group?”. By making use of the response-function parameterization, we can obtain a bound that does not depend on the specific parameters of the SCM, yielding a “worst-case” lower bound and a “best-case” upper bound on the probability that the patient would then be cured.

**Proposition 8.5.1** (Bounding counterfactual probabilities). *For a simple SCM  $M$  with a single binary exogenous input variable  $X \in \{0, 1\}$  and a binary endogenous variable*

$Y \in \{0, 1\}$  (and perhaps additional endogenous variables, and with an arbitrary number of exogenous variables taking values in arbitrary standard measurable spaces), the counterfactual probability

$$P_{M^{\text{twin}}}((Y')^{\text{do}(x'=1)} = 1 \mid Y^{\text{do}(x=0)} = 0)$$

(with  $Y^{\text{do}(x=0)}$  the factual outcome, and  $(Y')^{\text{do}(x'=1)}$  the counterfactual outcome) is bounded by

$$\frac{q_{0|0} - \min(q_{0|0}, q_{0|1})}{q_{0|0}} \leq P_{M^{\text{twin}}}((Y')^{\text{do}(x'=1)} = 1 \mid Y^{\text{do}(x=0)} = 0) \leq \frac{\min(q_{0|0}, q_{1|1})}{q_{0|0}},$$

where  $q_{y|x} := P_M(Y = y \mid \text{do}(X = x))$ .

*Proof.* Denote the 3-dimensional probability simplex by  $\Delta_3 = \{r = (r_{00}, r_{01}, r_{10}, r_{11}) \in [0, 1]^4 : r_{00} + r_{01} + r_{10} + r_{11} = 1\}$ . We know from Corollary 8.4.3 that without loss of generality, we may assume that the SCM has only a single binary endogenous variable  $Y$  and an exogenous random variable taking values in  $\{0, 1\}^{\{0,1\}}$ . That SCM must then lie in the family  $\{M_\rho : \rho \in \Delta_3\}$ , where the SCM  $M_\rho$  with parameter  $\rho$  has binary exogenous input variable  $X$ , binary endogenous variable  $Y$ , a single latent exogenous random variable  $W \in \{f_{00}, f_{01}, f_{10}, f_{11}\}$ , exogenous distribution  $P(W = f_w) = \rho_w$ , and structural equation

$$Y^{\text{do}(x)} = W(x),$$

where we defined response functions  $f_{00}, f_{01}, f_{10}, f_{11} : \{0, 1\} \rightarrow \{0, 1\}$  by:

$$\begin{aligned} f_{00} &: 0 \mapsto 0, 1 \mapsto 0; \\ f_{01} &: 0 \mapsto 0, 1 \mapsto 1; \\ f_{10} &: 0 \mapsto 1, 1 \mapsto 0; \\ f_{11} &: 0 \mapsto 1, 1 \mapsto 1. \end{aligned}$$

We will derive a bound on the counterfactual probability

$$P_{M_\rho^{\text{twin}}}((Y')^{\text{do}(x'=1)} = 1 \mid Y^{\text{do}(x=0)} = 0).$$

We first update the distribution of  $W$  with the observed outcome:

$$P_{M_\rho^{\text{twin}}}(W = w \mid Y^{\text{do}(x=0)} = 0) = \begin{cases} \frac{\rho_{00}}{\rho_{00} + \rho_{01}} & w = f_{00}, \\ \frac{\rho_{01}}{\rho_{00} + \rho_{01}} & w = f_{01}, \\ 0 & w = f_{10}, \\ 0 & w = f_{11}. \end{cases}$$

Because of the counterfactual equivalence of  $M_\rho$  and  $M$  w.r.t.  $Y$ , we have that  $q_{y|x} := P_M(Y = y \mid \text{do}(X = x)) = P_{M_\rho^{\text{twin}}}(Y = y \mid \text{do}(X = x))$ . This Markov kernel is given explicitly by:

$x$	$y$	$q_{y x}$
0	0	$\rho_{00} + \rho_{01}$
0	1	$\rho_{10} + \rho_{11}$
1	0	$\rho_{00} + \rho_{10}$
1	1	$\rho_{01} + \rho_{11}$

For our particular counterfactual probability of interest, we have the equality

$$P_{M_\rho^{\text{twin}}}((Y')^{\text{do}(x'=1)} = 1 | Y^{\text{do}(x=0)} = 0) = \frac{\rho_{01}}{q_{0|0}}.$$

From the non-negativity of the  $\rho$ 's and the table above, we can derive the bound

$$q_{0|0} - \min(q_{0|0}, q_{0|1}) \leq \rho_{01} \leq \min(q_{0|0}, q_{1|1})$$

and hence

$$\frac{q_{0|0} - \min(q_{0|0}, q_{0|1})}{q_{0|0}} \leq P_{M_\rho^{\text{twin}}}((Y')^{\text{do}(x'=1)} = 1 | Y^{\text{do}(x=0)} = 0) \leq \frac{\min(q_{0|0}, q_{1|1})}{q_{0|0}}.$$

Since  $M_\rho$  is counterfactually equivalent to the original SCM  $M$  w.r.t.  $Y$ , the same bound also holds for  $M$  instead of  $M_\rho$ . □

As an illustration, suppose that  $q_{0|0} \approx 1$  and  $q_{1|1} \approx 1$ . Then the bound tells us that  $P_{M^{\text{twin}}}((Y')^{\text{do}(x'=1)} = 1 | Y^{\text{do}(x=0)} = 0) \approx 1$  as well. Thus, for almost-deterministic relations, we can tightly bound this counterfactual probability.

## 9. Causal Discovery

So far, we always assumed that an SCM was fully specified, and derived theory to draw conclusions from the given SCM. For example, the do-calculus provides precise relationships between certain Markov kernels induced by the SCM, which enables us to perform *causal reasoning*.

However, we often do not have sufficient information regarding the system that we are modeling to completely specify an SCM. For example, we may only know what the observed variables are, but not what the graph of the SCM is, let alone know the latent spaces, exogenous distribution and exact causal mechanisms. Can we still perform causal reasoning with such incompletely specified models? The answer turns out to be affirmative, if one is willing to make certain assumptions (that—unfortunately, but perhaps unavoidably—are typically untestable).

In the rest of this chapter we will focus on the question of how to deduce partial knowledge about the graph from given Markov kernels. This is often called *causal discovery*. In the next chapter, we will go one step further, and replace the deduction of graphical properties from Markov kernels by the inference of graphical properties from data, i.e., we replace Markov kernels by finite samples. This will lead to statistical considerations. For example, one might study the properties of different estimators of causal effects. Estimating the causal effect of some variables on others is often called *causal inference* (although “inference” is often interpreted much broader as drawing conclusions from data and prior beliefs).

In this chapter, we will make use of the simple SCM formalism. Similar (more restrictive) results can be obtained in the L-CBN formalism.

### 9.1. Detecting Causal Relations

In this chapter, we will frequently be making use of the following short-hand notation.

**Definition 9.1.1.** *Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V$ . For  $O \subseteq V$ , we write  $G_{O|J} := G^{\setminus(V \setminus O)}$  for the graph obtained by marginalizing all output nodes except for those in  $O$ .*

We start by formalizing Definition 1.2.3 for simple SCMs.

**Definition 9.1.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G = G(M)$ . If  $a \in \text{Anc}^G(b) \setminus \{b\}$  for  $a, b \in J \cup V \cup W$  we say that  $a$  is a cause of  $b$  according to  $M$ .*

**Remark 9.1.3.** *With Remark 3.2.15 one sees that this holds if and only if  $a \rightarrow b$  is present in  $(G_{\text{do}(a)})_{b|J}$ .*

With the help of the do-calculus, we can now tie this notion to practical procedures to deduce the presence of causal relations from the (intervened) Markov kernels of the SCM. The following proposition expresses that only if  $a$  causes  $b$  according to  $M$ , a hard intervention on  $a$  can change the distribution of  $b$ . This formalizes our intuitive notion of what it means for a variable to cause another variable.



**Proposition 9.1.4.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let  $a \in V \cup W \cup J$  and  $b \in V$ .*

- For  $a \in V \cup W$ :

$$a \notin \text{Anc}^{G(M)}(b) \implies P_M(X_b | \text{do}(X_J), \text{do}(X_a)) = P_M(X_b | \text{do}(X_J));$$

- for  $a \in J$ :

$$a \notin \text{Anc}^{G(M)}(b) \implies P_M(X_b | \text{do}(X_J)) = P_M(X_b | \text{do}(X_{J \setminus \{a\}}), \underline{\text{do}}(X_a)).$$

*Note: both cases can also be combined into a single statement:*

$$a \notin \text{Anc}^{G(M)}(b) \implies P_M(X_b | \text{do}(X_{J \cup \{a\}})) = P_M(X_b | \text{do}(X_{J \setminus \{a\}})).$$

*Proof.* Denote  $G = G(M)$ . Assume that  $a \notin \text{Anc}^G(b)$ .

Let us first consider the case  $a \in V \cup W$ . We will show that

$$b \underset{G_{\text{do}(I_a)}}{\overset{\sigma}{\perp}} I_a | J \setminus \{a\},$$

where  $I_a$  is an intervention node with target  $a$  corresponding to an exogenous input variable modeling a hard intervention on  $a$ . Assume on the contrary that there exists a walk in  $G_{\text{do}(I_a)}$  between  $b$  and  $I_a \cup J$  that is  $\sigma$ -open given  $J \setminus \{a\}$ , and for which all colliders lie in  $J \setminus \{a\}$ . It cannot contain a node from  $J \setminus \{a\}$ , since that would either be an end node ( $\sigma$ -blocking the walk) or a non-collider node pointing only to nodes in another strongly connected component ( $\sigma$ -blocking the walk). Therefore it must be of the form  $I_a \rightarrow a \rightarrow \dots b$ . If it were a directed walk from  $I_a$  all the way to  $b$ , then we would get a contradiction with  $a \notin \text{Anc}^G(b)$ . Therefore, it must contain a collider. This collider must be in  $J \setminus \{a\}$ , which is a contradiction (since no input node can be a collider on a walk). We now apply rule 3 of the do-calculus (Theorem 7.3.1) to obtain:

$$P_M(X_b | \text{do}(X_J), \text{do}(X_a)) = P_M(X_b | \text{do}(X_J)).$$

Now consider the case  $a \in J$ . In that case we do not add an intervention node. We will show that

$$b \underset{G}{\overset{\sigma}{\perp}} a | J \setminus \{a\}.$$

Assume on the contrary that there exists a walk in  $G$  between  $b$  and  $J$  that is  $\sigma$ -open given  $J \setminus \{a\}$ , and for which all colliders lie in  $J \setminus \{a\}$ . It cannot contain a node from  $J \setminus \{a\}$ , since that would either be an end node ( $\sigma$ -blocking the walk) or a non-collider node pointing only to nodes in another strongly connected component ( $\sigma$ -blocking the walk). Therefore it must be of the form  $a \rightarrow \dots b$ . If it were a directed walk from  $a$  all the way to  $b$ , then we would get a contradiction with  $a \notin \text{Anc}^G(b)$ . Therefore, it must contain a collider. This collider must be in  $J \setminus \{a\}$ , which is a contradiction (since no input node can be a collider on a walk). We now apply rule 4 of the do-calculus (Theorem 7.3.1) to obtain:

$$P_M(X_b | \text{do}(X_J)) = P_M(X_b | \text{do}(X_{J \setminus \{a\}}), \underline{\text{do}}(X_a)).$$

□

This leads to a practical way of detecting the presence of a causal relation.

**Corollary 9.1.5.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let  $b \in V$ .*

*For  $a \in V \cup W \cup J$  with  $a \neq b$ , if there exist values  $x_{J \setminus \{a\}} \in \mathcal{X}_{J \setminus \{a\}}$ ,  $x_a, x'_a \in \mathcal{X}_a$  such that*

$$\begin{aligned} & P_M(X_b \mid \text{do}(X_a = x_a), \text{do}(X_{J \setminus \{a\}} = x_{J \setminus \{a\}})) \\ & \neq P_M(X_b \mid \text{do}(X_a = x'_a), \text{do}(X_{J \setminus \{a\}} = x_{J \setminus \{a\}})), \end{aligned}$$

*then  $a$  is a cause of  $b$  according to  $M$ , that is,  $a \in \text{Anc}^{G(M)}(b)$ .*

*Also, for  $a \in V \cup W$  with  $a \neq b$ , if there exist values  $x_J \in \mathcal{X}_J$ ,  $x_a \in \mathcal{X}_a$  such that*

$$P_M(X_b \mid \text{do}(X_a = x_a), \text{do}(X_J = x_J)) \neq P_M(X_b \mid \text{do}(X_J = x_J)),$$

*then  $a$  is a cause of  $b$  according to  $M$ .*

This condition is sufficient, but not necessary. This formalizes the main principle of how we can learn about causal relations in the world: by actively changing some part of the world (choosing the intervention values independently) and observing the response of other parts of the world. The independence assumption is key to distinguish mere correlation from causation.<sup>44</sup>

## 9.2. Detecting Direct Causal Relations

Another popular notion is that of direct causation. One should keep in mind that this is always relative to some set of variables. In particular, this property is not necessarily preserved under marginalization.

**Definition 9.2.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. If  $a \in \text{Pa}^{G(M)}(b)$  for  $a, b \in J \cup V \cup W$  with  $a \neq b$  we say that  $a$  is a direct cause of  $b$  w.r.t.  $V \cup W \cup J$  according to  $M$ .*

**Proposition 9.2.2.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. For  $a, b \in J \cup V \cup W$ ,  $a$  is a direct cause of  $b$  w.r.t.  $J \cup V \cup W$  according to  $M$  if and only if  $a$  is a cause of  $b$  according to  $M_{[b]}$ .*

Applying Proposition 9.1.4 to the intervened SCM  $M_{[b]}$  gives similar conditions to identify the presence of direct causal relations.

---

<sup>44</sup>As a less mathematical and more philosophical footnote: it is interesting to speculate about how this relates to the notion of a free will. If an agent is not convinced that it chose the intervention values independently of other past aspects of the world, it cannot validly perform this causal reasoning step. An agent without a free will to choose these values could therefore never conclude that its actions have a causal effect on the world, as it could also just be a puppet steered by higher powers, and any dependence it observes between its actions and aspects of the world could also be ascribed to confounding. So perhaps that is why evolution equipped us with the impression that we have a free will.

**Corollary 9.2.3.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let  $b \in V$ .*

*For  $a \in V \cup W \cup J$  with  $a \neq b$ , if there exist values  $x_{J \setminus \{a\}} \in \mathcal{X}_{J \setminus \{a\}}$ ,  $x_{V \setminus \{a,b\}} \in \mathcal{X}_{V \setminus \{a,b\}}$ ,  $x_a, x'_a \in \mathcal{X}_a$  such that*

$$\begin{aligned} & P_M(X_b \mid \text{do}(X_a = x_a), \text{do}(X_{V \setminus \{a,b\}} = x_{V \setminus \{a,b\}}), \text{do}(X_{J \setminus \{a\}} = x_{J \setminus \{a\}})) \\ & \neq P_M(X_b \mid \text{do}(X_a = x'_a), \text{do}(X_{V \setminus \{a,b\}} = x_{V \setminus \{a,b\}}), \text{do}(X_{J \setminus \{a\}} = x_{J \setminus \{a\}})), \end{aligned}$$

*then  $a$  is a direct cause of  $b$  w.r.t.  $V \cup J \cup W$  according to  $M$ .*

*Also, for  $a \in V \cup W$  with  $a \neq b$ , if there exist values  $x_J \in \mathcal{X}_J$ ,  $x_{V \setminus \{a,b\}} \in \mathcal{X}_{V \setminus \{a,b\}}$ ,  $x_a \in \mathcal{X}_a$  such that*

$$\begin{aligned} & P_M(X_b \mid \text{do}(X_a = x_a), \text{do}(X_{V \setminus \{a,b\}} = x_{V \setminus \{a,b\}}), \text{do}(X_J = x_J)) \\ & \neq P_M(X_b \mid \text{do}(X_{V \setminus \{a,b\}} = x_{V \setminus \{a,b\}}), \text{do}(X_J = x_J)), \end{aligned}$$

*then  $a$  is a direct cause of  $b$  w.r.t.  $V \cup J \cup W$  according to  $M$ .*

*Thus in both cases,  $a \in \text{Pa}^{G(M)}(b)$ .*

*Proof.* Apply Corollary 9.1.5 to the intervened SCM

$$M_{[b]} = (J \cup V \setminus \{b\}, \{b\}, W, \mathcal{X}, P, f_{\{b\}}),$$

and note that  $a \in \text{Anc}^{G(M_{[b]})}(b) \implies a \in \text{Pa}^{G(M)}(b) \cup \{b\}$ . □

Note further that this method to identify a direct causal effect may not be very practical, as it requires intervening on *all* endogenous and exogenous input variables (except  $b$ ) simultaneously.

### 9.3. Detecting Common Causes

The notion of having a common cause is formalized as follows.

**Definition 9.3.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G(M) = (J, V \cup W, E)$ . If there exists a bifurcation with source  $c \in V \cup W \cup J$  in  $G(M)$  between  $a, b \in V$ , we say that  $c$  is a common cause of  $a$  and  $b$  according to  $M$ .*

**Remark 9.3.2.** *With Remark 3.2.15 we get that this holds if and only if  $a \leftrightarrow b$  or  $a \leftarrow c \rightarrow b$  with  $c \in J$  is present in  $G_{\{a,b\}|J}$ .*

To test whether some variable is a common cause of two other variables, we can simply apply Corollary 9.1.5 for detecting causal relations (after appropriate interventions), making use of Proposition 3.2.4 to reexpress the existence of a bifurcation with source in terms of certain ancestral relations after performing certain interventions. We will not write this out in detail here.

Importantly: without common causes (and without a reverse causal relation), there can be no bias due to confounding.<sup>45</sup>

**Proposition 9.3.3.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let  $a \neq b \in V$ . If  $b \notin \text{Anc}^{G(M)}(a)$  and  $a$  and  $b$  do not have a common cause in  $V \cup W$  according to  $M$ , then  $a$  and  $b$  do not have confounding bias. That is,*

$$P_M(X_b \mid \text{do}(X_a), \text{do}(X_J)) = P_M(X_b \mid X_a, \text{do}(X_J)) \quad \mu_a\text{-a.s.}$$

for any reference measure  $\mu_a$  on  $\mathcal{X}_a$  with  $\mu_a \ll P_M(X_a \mid \text{do}(X_J)) \ll \mu_a$ .

*Proof.* Denote  $G = G(M)$ . We will show that the assumptions imply  $b \perp_{G_{\text{do}(I_a)}}^\sigma I_a \mid \{a\} \cup J$ .

Suppose on the contrary that there exists a walk in  $G_{\text{do}(I_a)}$  between  $b$  and  $I_a \cup J$  in  $G_{\text{do}(I_a)}$  that is  $(J \cup \{a\})$ - $\sigma$ -open. It cannot contain nodes from  $J$ , as such a node would either be an end node of the walk ( $\sigma$ -blocking it) or a non-collider node pointing only to nodes in another strongly connected component of  $G_{\text{do}(I_a)}$  ( $\sigma$ -blocking the walk). Hence it must be of the form  $I_a \rightarrow a \cdots b$ . Then there exists a walk  $I_a \rightarrow a \cdots b$  in  $G_{\text{do}(I_a)}$  that is  $(J \cup \{a\})$ - $\sigma$ -open and contains at least one collider. Indeed, suppose we have such a walk without a collider. Then it must be a directed walk  $I_a \rightarrow a \rightarrow \cdots \rightarrow d \rightarrow \cdots \rightarrow b$  where  $a \rightarrow \cdots \rightarrow d$  is the longest subwalk that spans a single strongly connected component of  $G_{\text{do}(I_a)}$  (equivalently: of  $G$ ). If  $d = a$ , the walk would not be  $(J \cup \{a\})$ - $\sigma$ -open. If  $d = b$ , we would get a contradiction with the assumption  $b \notin \text{Anc}^G(a)$ . Therefore,  $d$  must point to a node in another strongly connected component. We can now replace the subwalk  $a \rightarrow \cdots \rightarrow d$  by a directed walk through the same strongly connected component in the other direction,  $a \leftarrow \cdots \leftarrow d$ . In this way we obtain the walk  $I_a \rightarrow a \leftarrow \cdots \leftarrow d \rightarrow \cdots b$  which is also  $(J \cup \{a\})$ - $\sigma$ -open, and contains  $a$  as a collider.

Therefore, there must exist a walk  $I_a \rightarrow a \cdots b$  in  $G_{\text{do}(I_a)}$  that is  $(J \cup \{a\})$ - $\sigma$ -open and contains a collider. Each collider on the walk must be  $a$ . There must exist such a walk of minimal length, which can contain only a single collider. That walk must be of the form  $I_a \rightarrow a \leftarrow \cdots b$ , where the part between  $a$  and  $b$  contains no colliders and is not a directed walk from  $b$  to  $a$ . Hence it must be of the form  $I_a \rightarrow a \leftarrow \cdots \leftarrow c \rightarrow \cdots \rightarrow b$ , with  $c \notin J$ , and where  $b$  does not appear in the subwalk between  $a$  and  $c$ , and  $a$  does not appear in the subwalk between  $c$  and  $b$ . Contradiction. Hence  $b \perp_{G_{\text{do}(I_a)}}^\sigma I_a \mid \{a\} \cup J$ .

<sup>45</sup>Another bias is still possible: selection bias. This may happen if a subset of the samples (statistical units, individuals, ...) do not end up in the data set because of some filtering process. In other words, even if

$$P_M(X_b \mid \text{do}(X_a), \text{do}(X_J)) = P_M(X_b \mid X_a, \text{do}(X_J)) \quad \mu_a\text{-a.s.}$$

we need not have

$$P_M(X_b \mid S \in \xi_S, \text{do}(X_a), \text{do}(X_J)) = P_M(X_b \mid S \in \xi_S, X_a, \text{do}(X_J)) \quad \mu_a\text{-a.s.}$$

for some  $S \subseteq V \cup W$  and measurable subset  $\xi_S \subseteq \mathcal{X}_S$ . We plan to add more discussion on this to a future version of these lecture notes.

Invoking rule 2 of the do-calculus (Theorem 7.3.1) then gives

$$P_M(X_b | \text{do}(X_a), \text{do}(X_J)) = P_M(X_b | X_a, \text{do}(X_J)) \quad \mu_a\text{-a.s.}$$

provided that  $\mu_a \ll P_M(X_a | \text{do}(X_J)) \ll \mu_a$ .  $\square$

This leads to the following criterion to detect the presence of a common cause:

**Corollary 9.3.4.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let  $a, b \in V$  with  $a \neq b$ . Let  $\mu_a$  be a reference measure on  $\mathcal{X}_a$  with  $\mu_a \ll P_M(X_a | \text{do}(X_J)) \ll \mu_a$ . If  $b$  is not a cause of  $a$  according to  $M$ , and the following does not hold:*

$$P_M(X_b | \text{do}(X_a), \text{do}(X_J)) = P_M(X_b | X_a, \text{do}(X_J)) \quad \mu_a\text{-a.s.},$$

then  $a$  and  $b$  have a common cause in  $V \cup W$  according to  $M$ .

Furthermore, if for some  $c \in J$ , there exist values  $x_{J \setminus \{c\}} \in \mathcal{X}_{J \setminus \{c\}}$ ,  $x_c, x'_c \in \mathcal{X}_c$  such that

$$\begin{aligned} & P_M(X_a | \text{do}(X_c = x_c), \text{do}(X_{J \setminus \{c\}} = x_{J \setminus \{c\}})) \\ & \neq P_M(X_a | \text{do}(X_c = x'_c), \text{do}(X_{J \setminus \{c\}} = x_{J \setminus \{c\}})), \end{aligned}$$

and there exist values  $x_{J \setminus \{c\}} \in \mathcal{X}_{J \setminus \{c\}}$ ,  $x_c, x'_c \in \mathcal{X}_c$  such that

$$\begin{aligned} & P_M(X_b | \text{do}(X_c = x_c), \text{do}(X_{J \setminus \{c\}} = x_{J \setminus \{c\}})) \\ & \neq P_M(X_b | \text{do}(X_c = x'_c), \text{do}(X_{J \setminus \{c\}} = x_{J \setminus \{c\}})), \end{aligned}$$

then  $c$  is a common cause of  $a$  and  $b$  according to  $M$ .

*Proof.* The first part is the contra-positive of Proposition 9.3.3. The second part applies Corollary 9.1.5 to obtain sufficient conditions for the existence of  $c \in J$  that causes both  $a$  and  $b$  according to  $M$ .  $\square$

This condition is sufficient, but not necessary. To apply it, we need to know already that  $b$  does not cause  $a$  according to  $M$ . By swapping the roles of  $a$  and  $b$ , we can also use this if we know that  $a$  does not cause  $b$ .<sup>46</sup> Furthermore, note that we have implicitly assumed here the absence of selection bias.

In the potential outcome literature, one encounters other criteria for unconfoundedness that are formulated in terms of counterfactuals.

**Proposition 9.3.5.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM. Let  $a \neq b \in V$ . Consider the twin SCM  $M^{\text{twin}}$ . If  $b$  does not cause  $a$  according to  $M$ , and  $a$  and  $b$  have no common cause according to  $M$ , then*

$$\forall x_{a'} \in \mathcal{X}_a : \quad X_a \perp\!\!\!\perp_{P_{M^{\text{twin}}_{\text{do}(X_{a'}=x_{a'}})}} X_b | X_J. \quad (48)$$

<sup>46</sup>How to detect common causes if  $a$  and  $b$  are part of a causal cycle is an open research problem.

*Proof.* Let  $x_{a'} \in \mathcal{X}_a$ . Suppose there exists a  $\sigma$ -open walk in the graph  $G(M_{\text{do}(X_{a'}=x_{a'})}^{\text{twin}}) = G(M^{\text{twin}})_{\text{do}(a')}$  between  $a$  and  $b' \cup J$ . Then there must be such a walk of minimal length. It cannot contain nodes from  $J$ , as such a node would either be an end node of the walk ( $\sigma$ -blocking it) or a non-collider node pointing only to nodes in another strongly connected component ( $\sigma$ -blocking the walk). Hence it must be between  $a$  and  $b'$ . It cannot contain any collider. It cannot be a directed walk because it has to pass through an exogenous random node in  $W$ , but those nodes have no incoming edges. Therefore it must be a walk of the form  $a \leftarrow \dots \leftarrow c \rightarrow \dots \rightarrow b'$ . But then  $c \in W$ , because only nodes in  $W$  can be ancestors of endogenous nodes in different “worlds”. Note that the subwalk  $a \leftarrow \dots \leftarrow c$  must consist of nodes in  $V \cup \{c\}$  and cannot contain  $b$ , otherwise we have a contradiction with the assumption that  $b$  does not cause  $a$  according to  $M$ . Also, the subwalk  $c \rightarrow \dots \rightarrow b'$  must consist of nodes in  $V' \cup \{c\}$  and cannot contain  $a'$ , as  $a'$  has no incoming edges in  $G(M^{\text{twin}})_{\text{do}(a')}$ . By “removing the primes” from this subwalk, we obtain a walk of the form  $a \leftarrow \dots \leftarrow c \rightarrow \dots \rightarrow b$  in  $G(M)$ , which is seen to be a bifurcation with source  $c$ . But then  $a$  and  $b$  would have a common cause according to  $M$ , contradicting the assumptions. Hence

$$a \perp_{G(M^{\text{twin}})_{\text{do}(a')}} b' \mid J.$$

By the global Markov property applied to  $M_{\text{do}(X_{a'}=x_{a'})}^{\text{twin}}$  we conclude that

$$X_a \perp_{P_{M^{\text{twin}}_{\text{do}(X_{a'}=x_{a'})}}} X_{b'} \mid X_J.$$

□

Equation (48), typically written in terms of potential outcomes as

$$\forall x_J = x_{J'} \in \mathcal{X}_J, \forall x_{a'} \in \mathcal{X}_a : \quad X_a^{\text{do}(x_J)} \perp\!\!\!\perp X_{b'}^{\text{do}(x_{J'}, x_{a'})},$$

and commonly encountered for the special case  $J = \emptyset$ ,

$$\forall x_{a'} \in \mathcal{X}_a : \quad X_a \perp\!\!\!\perp X_{b'}^{\text{do}(x_{a'})},$$

usually written as

$$\forall x_{a'} \in \mathcal{X}_a : \quad X_a \perp\!\!\!\perp X_{b'}(x_{a'}),$$

is referred to as “exchangeability” in the potential outcome framework, and expresses the assumption of “no confounding” when the task is to estimate the causal effect of  $a$  on  $b$ .

## 9.4. Abstraction through Marginalization

We often deal with the situation that only part of the variables are observed or observable, and others are latent.

**Notation 9.4.1.** For a simple SCM  $M$ , let  $O \subseteq V \cup W$  be the set of observable variables, where we will always assume  $J$  to be observable as well. We refer to the marginalized graph

$$G_{O|J}(M) := (G(M)) \setminus ((V \cup W) \setminus O)$$

as the observable graph of  $M$ . It has input nodes  $J$  and output nodes  $O$ . The corresponding observable Markov kernel is denoted

$$P_M(X_O \mid \text{do}(X_J)),$$

and the observable intervened Markov kernel resulting from a hard intervention targeting  $T \subseteq O$  is

$$P_M(X_{O \setminus T} \mid \text{do}(X_{J \cup T})) := P_{M_{\text{do}(T)}}(X_{O \setminus T} \mid \text{do}(X_{J \cup T})).$$

Remember that Remark 3.2.15 stated that graphical marginalization preserves ancestral relations and bifurcations, while Lemma 3.3.8 stated that graphical marginalization preserves d-separation statements and  $\sigma$ -separation statements concerning observable variables. These properties have powerful consequences:

**Remark 9.4.2.** For many causal reasoning purposes, it suffices to know the observable graph  $G_{O|J}(M)$  rather than  $G(M)$ . In particular, the observable graph contains all information concerning:

1. which pairs of observable variables are causally related according to  $M$  (c.f. Definition 9.1.2);
2. which pairs of observable variables have a common cause according to  $M$  (c.f. Definition 9.3.1);
3. the d-separation and  $\sigma$ -separation statements between observable variables according to  $M$  (c.f. Definition 3.3.5).

In addition, the probabilistic marginalization preserves the notion of conditional independence. That is, for  $A, B, C \subseteq J \cup O$ , we have

$$X_A \perp\!\!\!\perp_{P_M(X_V, X_W | X_J)} X_B \mid X_C \iff X_A \perp\!\!\!\perp_{P_M(X_O | X_J)} X_B \mid X_C.$$

This means that for applying the global Markov property on observable variables (and hence the do-calculus and adjustment criteria), it suffices to know the observable graph  $G_{O|J}(M)$  and the observable (intervened) Markov kernels rather than the full graph  $G(M)$  and the full Markov kernels associated to  $M$ . **In other words: these properties taken together are very powerful, as they allow us to abstract away irrelevant details when performing causal reasoning.**

The following exercise illustrates this.

**Exercise 9.4.3.** *Reichenbach’s Principle of Common Cause states that if two events are dependent, then one must cause the other or the events must have a common cause (or any combination of these three possibilities).*

*We can make this precise for simple SCMs in the following way. Assume that  $M$  is a simple SCM with two observed endogenous variables  $X, Y$  (and possibly other latent variables as well, but no exogenous input nodes). Prove that  $X \not\perp_{P_M(X,Y)} Y$  implies that  $X \rightarrow Y$ ,  $X \leftarrow Y$  or  $X \leftrightarrow Y$  in  $G_{\{X,Y\}}(M)$ .*

*Thus, either  $X$  causes  $Y$  according to  $M$ , or  $Y$  causes  $X$  according to  $M$ , or  $X$  and  $Y$  have a common cause according to  $M$ .*

## 9.5. Randomized Controlled Trials

The notion of randomized controlled trials (also known as A/B-testing in engineering), is centuries old. It was already proposed by the Flemish physician Jan Baptista van Helmont [vH48] in 1648. As of today, it still provides the ‘gold standard’ for discovering causal relations and for the estimation of causal effects.

The experimental procedure is as follows. Consider two variables, “treatment”  $C$  and “outcome”  $X$ . In the simplest setting, one considers a binary treatment variable, where  $C = 1$  corresponds to “treat with drug” and  $C = 0$  corresponds to “treat with placebo” in a medical setting, or with “arm A” and “arm B” in an engineering setting. For example, the drug could be aspirin, and outcome could be the severity of headache perceived two hours later. Patients are split into two groups, the treatment and the control group, by means of a coin flip that assigns a value of  $C$  to every patient.<sup>47</sup> Patients are treated depending on the assigned value of  $C$ , i.e., patients in the treatment group are treated with the drug and patients in the control group are treated with a placebo. Some time after treatment, the outcome  $X$  is measured for each patient. This yields a data set  $(C_n, X_n)_{n=1}^N$  with two measurements  $(C_n, X_n)$  for the  $n^{\text{th}}$  patient. If the distribution of outcome  $X$  significantly differs between the two groups, one concludes that treatment is a cause of outcome.

Let us formalize this in the causal modeling language of SCMs. Apart from that treatment may have a causal effect on outcome, there are likely many other factors that influence outcome. Some have been measured, others not. For obvious practical reasons, we are not going to explicitly model *all* of them. Formally, we will assume that an accurate causal model of the situation is provided by some (unknown) simple SCM with observed variables  $C$  and  $X$ , and possibly other latent variables. We will consider the outcome variable  $X$  as endogenous. But what type of variable should we consider the treatment variable  $C$  to be (which is not necessarily binary)? We have three possibilities: exogenous input, exogenous random, and endogenous. We will discuss each of these three possibilities in sequence.

Let us start by considering the treatment variable  $C$  as an exogenous input variable. We are interested in answering two questions. The first is “Does treatment cause out-

---

<sup>47</sup>Usually this is done in a double-blind way, so that neither the patient nor the doctor knows which group a patient has been assigned to.



come?”, where we interpret this question as that the hypothetical causal relation should hold according to the underlying SCM  $M$ . In terms of the observable graph  $G_{X|C}(M)$ , this is then equivalent to asking “Is  $C \rightarrow X$  in  $G_{X|C}(M)$ ?”. The second question is “What is the causal effect of treatment on outcome?”. We interpret this as asking for the Markov kernel  $P_M(X | \text{do}(C))$ .

**Proposition 9.5.1.** *Let  $M$  be a simple SCM with a single exogenous input variable  $C$  and an endogenous variable  $X$ , both of which are observed (and possibly other latent variables as well). A dependence*

$$X \not\perp_{P_M(X|\text{do}(C))} C \quad (49)$$

*implies that  $C$  causes  $X$  according to  $M$ .*

*Proof.* This follows immediately from Proposition 9.1.4. □

Alternatively, we can consider the treatment variable as an exogenous *random* variable.

**Proposition 9.5.2.** *Let  $\bar{M}$  be a simple SCM with an exogenous random variable  $C$  and an endogenous variable  $X$ , both of which are observed (and possibly other latent variables as well). A dependence*

$$X \not\perp_{P_{\bar{M}}(X,C)} C \quad (50)$$

*implies that  $C$  causes  $X$  according to  $\bar{M}$ . The causal effect of  $C$  on  $X$  satisfies:*

$$P_{\bar{M}}(X | \text{do}(C)) = P_{\bar{M}}(X | C) \quad P_{\bar{M}}(C)\text{-a.s.} \quad (51)$$

*Proof.* Denote the observable graph as  $\bar{G} := G_{\{C,X\}}(\bar{M})$ . It has two nodes, and it either has no edge at all, in which case  $C$  does not cause  $X$  according to  $\bar{M}$ , or it has a single edge  $C \rightarrow X$ , in which case  $C$  causes  $X$  according to  $\bar{M}$ . By the Markov property (Corollary 7.2.1), if the edge  $C \rightarrow X$  were absent in  $\bar{G}$ , then  $X \perp_{P_{\bar{M}}(X,C)} C$ . In both cases, rule 2 of the causal do-calculus applied to  $\bar{G}$  yields the identity (51). □

A third option is to consider the treatment variable as *endogenous*. One situation in which this makes sense is so-called “imperfect compliance”. If trial subjects do not all comply with prescribed treatment, for whatever reasons, then we can no longer identify the coin flip outcome with treatment, (even though coin flip outcome may still be an important cause of treatment). In this modeling variant, we assume the existence of a simple SCM  $M$  with endogenous variables  $C, X$ , both of which are observed, and possibly other latent variables that provides an accurate model.  $C$  here retains the meaning of treatment (but is no longer necessarily identifiable with the coin flip result). Under additional assumptions regarding the exogeneity of the treatment variable, we again obtain a similar statement as before.

**Proposition 9.5.3.** *Let  $\tilde{M}$  be a simple SCM with two observed endogenous variables  $C, X$  (and possibly other latent variables as well). Under the following two assumptions:*

1.  $X$  does not cause  $C$  according to  $\tilde{M}$ , and
2.  $C$  and  $X$  do not have a common cause according to  $\tilde{M}$ ,

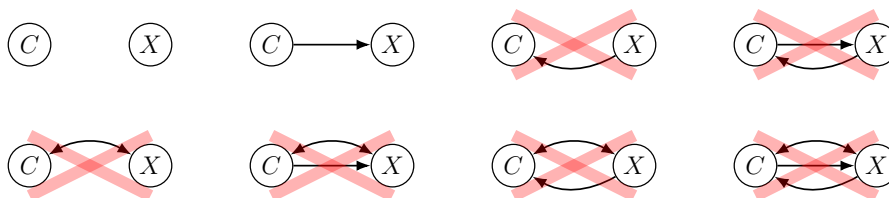
a dependence

$$X \not\perp_{P_{\tilde{M}}(X,C)} C \quad (52)$$

implies that  $C$  causes  $X$  according to  $\tilde{M}$ , and the causal effect of  $C$  on  $X$  satisfies:

$$P_{\tilde{M}}(X | \text{do}(C)) = P_{\tilde{M}}(X | C) \quad P_{\tilde{M}}(C)\text{-a.s.} \quad (53)$$

*Proof.* Denote the observable graph as  $\tilde{G} := G_{\{C,X\}}(\tilde{M})$ . The first assumption is equivalent to  $C \leftarrow X \notin \tilde{G}$ , and the second assumption is equivalent to  $C \leftrightarrow X \notin \tilde{G}$ . Hence, out of the eight possible graphs  $\tilde{G}$ , only two satisfy the assumptions:



By the Markov property, if the edge  $C \rightarrow X$  were absent in  $\tilde{G}$ , then  $X \perp_{P_{\tilde{M}}(X,C)} C$ . In both cases, rule 2 of the causal do-calculus applied to  $\tilde{G}$  yields the identity (53).  $\square$

Equation (49) is equivalent to the existence of values  $c, c' \in \mathcal{X}_C$  such that

$$P_M(X | \text{do}(C = c)) \neq P_M(X | \text{do}(C = c')).$$

Equation (50) is equivalent to the existence of values

$$P_{\tilde{M}}(X | C = c) \neq P_{\tilde{M}}(X | C = c')$$

for every version of  $P_{\tilde{M}}(X | C)$  (and something similar holds for equation (52)). These two statements are subtly different. We will see in the next chapter, that as long as  $C$  is discrete, they are actually not that different when testing these statements from a finite sample.

Apart from assuming that there exists a simple SCM that provides an accurate model, in all three cases, we made the following (implicit or explicit) causal assumptions regarding the treatment variable:

1. outcome  $X$  does not cause treatment  $C$ ;
2. outcome  $X$  and treatment  $C$  have no common cause, which implies that the values for the treatment variable are assigned independently of other (latent) factors that may influence the outcome.

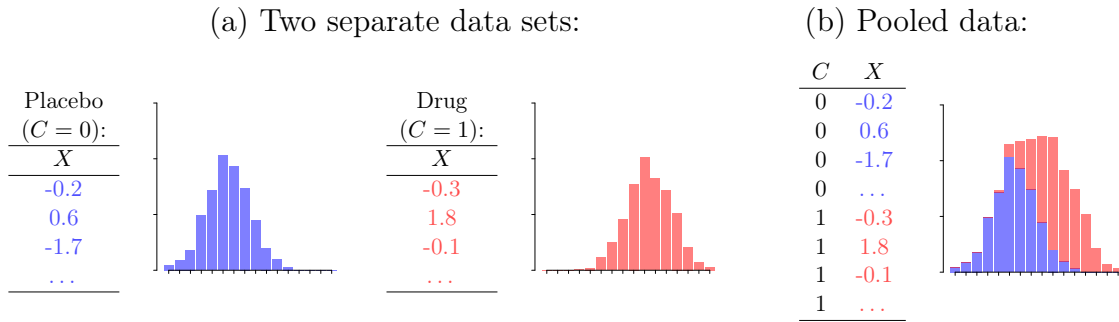


Figure 22: Illustration of the data from an example randomized controlled trial. When treatment  $C$  is randomized, the data can either be interpreted as (a) two separate data sets, one for the treatment and one for the control group, or (b) as a single data set including a context variable indicating treatment/control. Note that in this particular example,  $C$  is dependent on  $X$  in the pooled data (or equivalently, the distribution of  $X$  differs between contexts  $C = 0$  and  $C = 1$ ), which implies that  $C$  is a cause of  $X$ .

The first assumption is commonly deemed justified if the outcome is an event that occurs later in time than the treatment event. The second assumption is usually defended by appealing to randomization. Indeed, if treatment is decided solely by a proper coin flip, then it seems reasonable to assume that there is no irreducible common cause of coin flip and outcome (where “irreducible” means that it cannot be separated into statistically independent separate causes of both).

**Exercise 9.5.4.** *Think like a conspiracy theorist and imagine situations in which the first assumption is not valid. Do the same for the second assumption.*

Another implicit assumption we made is that the data was not subject to selection bias. In other words, no data is missing (except perhaps completely at random). For example, if patients with a negative outcome are more likely to report their outcome in a questionnaire than patients with a positive outcome, then this assumption may be violated.

We have shown (in three slightly different ways) that under these assumptions, if the distribution of the outcome  $X$  differs between the two groups of patients (“treatment group” with  $C = 1$  vs. “control group” with  $C = 0$ ), then treatment must be a cause of outcome, at least in this population of patients. Supposing that treatment is completely randomized, there are two conceptually different ways of testing this in the data, depending on whether we treat the data as a single pooled data set, or rather as two separate data sets (each one corresponding to a particular patient group), see also Figure 22. To test whether  $P(X | \text{do}(C = 0)) \neq P(X | \text{do}(C = 1))$ , we can test whether the distribution of  $X$  is statistically different in the two groups. This can be tested with a two-sample test, for example, a  $t$ -test or a Wilcoxon test. The other alternative is to consider the data as a single *pooled* data set. The question then becomes whether the conditional distribution of  $X$  given  $C = 0$  differs from the conditional distribution of

$X$  given  $C = 1$ , i.e., whether  $P(X | C = 0) \neq P(X | C = 1)$ . This can be done with a conditional independence test, for example, Fisher's exact test, or a partial correlation test. In the next lecture, we will look in more detail at possible tests.

In the end, the three ways of modeling the RCT are only slightly different. If one formally considers treatment as an exogenous input variable, but then also assumes that its values are randomly assigned, then the differences are purely cosmetic. However, there is one advantage that the exogenous input approach has over the other two: here we do not model *at all* how the values of treatment are chosen (except for the exogeneity assumptions). This allows more freedom in the experimental design and sampling scheme design. For example, one can decide ahead of the RCT that the sampling scheme should end up with an equal number of patients in both groups. In case treatment is assigned by flipping a coin for each patient, it is rather unlikely that we end up with exactly the same number of patients in both groups.

A fourth way to formalize the randomized controlled setting is by using potential outcomes. For a binary treatment variable, we introduce two random variables per patient:  $X_n^{\text{do}(c_n=1)}$  and  $X_n^{\text{do}(c_n=0)}$ , corresponding to the potential outcomes for the  $n$ 'th patient if we treat the patient, or not, respectively. Given the actual treatment  $C_n$ , we then define the actual outcome as  $X_n := X_n^{\text{do}(c_n=C_n)}$ . In practice, we only observe the actual outcome, and the other potential outcome remains latent.

## 9.6. Estimating average treatment effects

The task of estimating the causal effect of treatment on outcome is then often formulated as estimating the *average treatment effect (ATE)*

$$\tau := \mathbb{E}(X_n^{\text{do}(c_n=1)} - X_n^{\text{do}(c_n=0)}).$$

To do so, one assumes that  $C_n$  is randomized. This motivates the assumption that treatment and outcome are unconfounded, i.e., with Proposition 9.3.5:

$$X_n^{\text{do}(c_n=0)} \perp\!\!\!\perp C_n, \quad X_n^{\text{do}(c_n=1)} \perp\!\!\!\perp C_n.$$

One can then show that the difference-in-means estimator

$$\hat{\tau} := \frac{1}{|n : C_n = 1|} \sum_{\substack{n=1 \\ C_n=1}}^N X_n - \frac{1}{|n : C_n = 0|} \sum_{\substack{n=1 \\ C_n=0}}^N X_n$$

is an unbiased, consistent estimator of the ATE  $\tau$ . Curiously enough, while we can speak of the difference  $X_n^{\text{do}(c_n=1)} - X_n^{\text{do}(c_n=0)}$  as the individual treatment effect, this is a fundamentally unobservable quantity; however, the average treatment effect can be estimated from observed data. In the SCM setting, we can think of the potential outcomes as counterfactuals in a twin SCM. However, when assuming an underlying SCM, there is no need to go to the counterfactual level, as one can simply define the ATE as

$$\tau := \mathbb{E}_M(X | \text{do}(C = 1)) - \mathbb{E}_M(X | \text{do}(C = 0)).$$

In the presence of observed covariates  $Z$ , one often considers also the *conditional average treatment effect (CATE)*, which we can define as

$$\mathbb{E}_M(X \mid \text{do}(C = 1), Z) - \mathbb{E}_M(X \mid \text{do}(C = 0), Z).$$

when assuming an underlying SCM. There is a large body of literature that considers the question of studying the (asymptotic) efficiency of estimators of the (conditional) average treatment effect. For a nice account of this surprisingly non-trivial inference problem, see e.g. [Wag20].

## 9.7. Faithfulness

The converse statement of the global Markov property for simple SCMs (Corollary 7.2.1) and for causal Bayesian networks (Theorem 4.2.1) is called “faithfulness”.

**Definition 9.7.1.** *Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G(M)$  and Markov kernel  $P_M(X_V, X_W \mid \text{do}(X_J))$ .  $M$  is called  $\sigma$ -faithful if for all  $A, B, C \subseteq J \cup V \cup W$  (not necessarily disjoint):*

$$A \underset{G(M)}{\overset{\sigma}{\perp}} B \mid C \iff X_A \underset{P_M(X_V, X_W \mid \text{do}(X_J))}{\perp\!\!\!\perp} X_B \mid X_C \quad (54)$$

*It is called  $d$ -faithful if for all  $A, B, C \subseteq J \cup V \cup W$  (not necessarily disjoint):*

$$A \underset{G(M)}{\overset{d}{\perp}} B \mid C \iff X_A \underset{P_M(X_V, X_W \mid \text{do}(X_J))}{\perp\!\!\!\perp} X_B \mid X_C \quad (55)$$

*For a subset  $O \subseteq V \cup W$ , we say that  $M$  is  $\sigma$ -faithful w.r.t.  $O$  if (54) holds for all (not necessarily disjoint)  $A, B, C \subseteq J \cup O$ , and we define  $d$ -faithful w.r.t.  $O$  analogously.*

In words, a simple SCM is called  $\sigma$ -faithful ( $d$ -faithful) if each conditional independence in the induced Markov kernel is due to a  $\sigma$ -separation ( $d$ -separation).

**Remark 9.7.2.** *These notions behave properly under marginalization: a simple SCM  $M$  is  $\sigma/d$ -faithful w.r.t.  $O$  if and only if its marginalization  $M_O$  is  $\sigma/d$ -faithful.*

Faithfulness may fail for various reasons:

- Deterministic relationships may lead to additional conditional independences, but are not exploited by the Markov property;
- Effects may cancel out;
- If cycles are present, and (i) all variables are discrete, or (ii) interactions are linear;
- If cycles are present, and the system is “perfectly adapting”.

An example of a deterministic relationship leading to a faithfulness violation is the following.

**Example 9.7.3.** Take an SCM with three endogenous variables  $X, Y, Z$  and two exogenous random variables  $U, W$ , with structural equations

$$X = 5, \quad Y = X + U, \quad Z = X + W.$$

Then  $Y \not\perp Z$  but  $Y \perp\!\!\!\perp Z$ . This simple (even acyclic) SCM is not faithful due to  $X$  being constant.

The next example illustrates how canceling effect may lead to a faithfulness violation.

**Example 9.7.4.** Take an SCM with three endogenous variables  $X, Y, Z$  and three exogenous random variables  $W_X, W_Y, W_Z$ , with structural equations

$$X = W_X, \quad Y = X + W_Y, \quad Z = Y - X + W_Z.$$

Then  $X \not\perp Z$  but  $X \perp\!\!\!\perp Z$ .

One can show that in certain special cases, the global Markov property in terms of  $d$ -separation even holds for simple SCMs.

**Proposition 9.7.5.** Let  $M = (J, V, W, \mathcal{X}, P, f)$  be a simple SCM with graph  $G(M)$  and Markov kernel  $P_M(X_V, X_W \mid \text{do}(X_J))$ . If  $J = \emptyset$  and one of the three conditions applies:

1. all spaces  $\mathcal{X}_v$  with  $v \in V$  are discrete, or
2. the causal mechanism  $f$  is affine and the exogenous distribution has a density w.r.t. Lebesgue measure, or
3.  $M$  is acyclic,

then for all  $A, B, C \subseteq J \cup V \cup W$  (not necessarily disjoint):

$$A \underset{G(M)}{\overset{d}{\perp}} B \mid C \implies X_A \underset{P_M(X_V, X_W \mid \text{do}(X_J))}{\perp\!\!\!\perp} X_B \mid X_C.$$

The proofs are given in [FM17].

## 9.8. Local Causal Discovery

Although the most reliable way to discover causal relations and to estimate their effects is by means of a randomized controlled trial, it is not always possible or feasible to perform such an experiment. One alternative is provided by the Local Causal Discovery (LCD) algorithm [Coo97].

LCD is a *constraint-based* causal discovery algorithm which means that it discovers causal relations by combining the results of conditional independence tests on data. It can be used for the purely observational causal discovery setting where certain background knowledge is available that is weaker than that for the randomized controlled trial. In particular, no randomization is necessary.

The basic idea behind the LCD algorithm is the following result of [Coo97] (originally formulated for L-CBNs, but easily generalized to simple SCMs):



Figure 23: All possible observable graphs detected by LCD.

**Proposition 9.8.1.** *Let  $M$  be a simple SCM with observed endogenous variables  $O = \{1, 2, 3\} \subseteq V$  and no exogenous input variables ( $J = \emptyset$ ). Suppose that it is  $\sigma$ -faithful (w.r.t.  $O$ ). If  $X_2$  is not a cause of  $X_1$  according to  $M$ , the following conditional (in)dependencies<sup>48</sup> in the observational distribution  $P_M(X_1, X_2, X_3)$*

$$X_1 \not\perp\!\!\!\perp X_2, \quad X_2 \not\perp\!\!\!\perp X_3, \quad X_1 \perp\!\!\!\perp X_3 \mid X_2$$

imply that the observable graph  $G_O(M)$  must be one of the three DMGs in Figure 23. Hence,

1.  $X_3$  is not a cause of  $X_2$  according to  $M$ ;
2.  $X_2$  is a direct cause of  $X_3$  w.r.t  $\{1, 2, 3\}$  according to  $M$ ;
3.  $X_2$  and  $X_3$  do not have a common cause according to  $M$ ;
4. the causal effect of  $X_2$  on  $X_3$  according to  $M$  is given by:

$$P_M(X_3 \mid \text{do}(X_2)) = P_M(X_3 \mid X_2) \quad P_M(X_2)\text{-a.s.} \quad (56)$$

*Proof.* The proof proceeds by enumerating all (possibly cyclic) DMGs on three variables that the observable graph  $G_O(M)$  could be, and ruling out the ones that do not satisfy the assumptions. The assumption that  $X_2$  is not a cause of  $X_1$  implies that there is no directed edge  $X_2 \rightarrow X_1$  in the graph  $G_O(M)$ . If there were an edge between  $X_1$  and  $X_3$ ,  $X_1 \perp\!\!\!\perp X_3 \mid X_2$  would not hold (faithfulness). Also, since  $X_1 \not\perp\!\!\!\perp X_2$ ,  $X_1$  and  $X_2$  must be adjacent (Markov property). Similarly,  $X_2$  and  $X_3$  must be adjacent.  $X_2$  cannot be a collider on any walk between  $X_1$  and  $X_3$  (faithfulness). Since the only possible edges between  $X_1$  and  $X_2$  are  $X_1 \rightarrow X_2$  and  $X_1 \leftrightarrow X_2$  (both of which are into  $X_2$ ), this means that there must be a directed edge  $X_2 \rightarrow X_3$ , but there cannot be a bidirected edge  $X_2 \leftrightarrow X_3$  or directed edge  $X_2 \leftarrow X_3$ . In other words, the only three possible graphs are the ones in Figure 23. The causal do-calculus applied to  $G_O(M)$  yields (56).  $\square$

In one of the first applications of LCD, it was discovered that nausea causes vomiting [SC99]. The next example provides another successful application of LCD.

**Example 9.8.2.** *PIP2 and PIP3 are phospholipids that play an important role in human immune system cells. Figure 24 shows a scatter plot of PIP2 and PIP3 expression levels, measured in individual human immune system cells, after activation of certain protein signaling cascades in these cells [SPP<sup>+</sup>05]. The measurements have been performed*

<sup>48</sup>Henceforth, we will no longer always explicitly write the Markov kernel as a subscript to the conditional independence symbol if it is clear from the context which Markov kernel is meant.

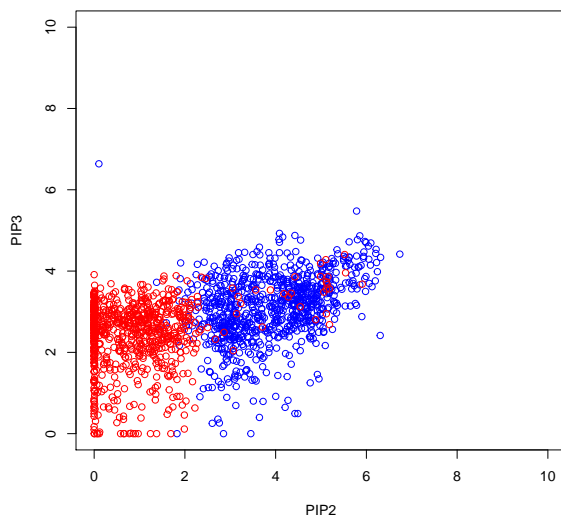


Figure 24: “Text-book” example of an LCD pattern in flow cytometry data of [SPP+05]. See Example 9.8.2 for details.

under two different experimental conditions: observational ( $C = 0$ , in blue), and after intervening by administering the chemical compound *Psitectorigenin* to the cells before measuring the  $PIP2$  and  $PIP3$  levels ( $C = 1$ , in red). The experimental protocol justifies the assumption that neither  $PIP2$  nor  $PIP3$  expression levels can cause the experimental condition (because these expression levels were measured after the experimental condition had been imposed on the cells). Also, assuming the cells to be properly randomized before split into the two groups (corresponding with the two experimental conditions) we can rule out the existence of a common cause of the experimental condition  $C$  and  $PIP2$ , and similarly of  $C$  and  $PIP3$ . The scatter plot suggests the following conditional independence:

$$PIP3 \perp\!\!\!\perp_{P(PIP2, PIP3 | do(C))} C | PIP2,$$

which can also be confirmed with statistical conditional independence tests. Therefore, we have an LCD pattern with  $X_1 = C$ ,  $X_2 = PIP2$ ,  $X_3 = PIP3$ , which allows us to infer that the  $PIP2$  expression level causes the  $PIP3$  expression level. Under the randomization assumption, we can even infer that *Psitectorigenin* exposure is a cause of  $PIP2$  expression levels. This is in line with the *Psitectorigenin* being known as an inhibitor of  $PIP2$ , reducing the quantity of  $PIP2$  in cells after exposure of the cells to this inhibitor.

A high-dimensional adaptation has also been shown to be successful in predicting the effects of gene knockout on gene expression levels from large-scale interventional yeast gene expression data [VM19].

In case more than three variables have been observed, one can run LCD on all triples of variables for which its assumptions apply. In that case, one should keep in mind that



a direct edge in a marginalized graph does not imply the presence of the directed edge in the original graph (only the presence of a directed path). In other words, with respect to a larger set of observed variables, the causal relations found by LCD are not necessarily direct.

In case of more than three observed variables, one can also replace the single variable  $X_2$  in the LCD algorithm by a subset of variables, a so-called *separating set*. This idea is exploited efficiently in case of many variables in the Invariant Causal Prediction algorithm [PBM16].

## 9.9. Y-structures

For both the randomized controlled trial and the LCD algorithm, we need prior knowledge: we need to know already that one of the variables is not a cause of another one. It turns out that in the absence of any such causal background knowledge, we can sometimes still deduce causal relationships from observed conditional independences. The simplest such example is given by the “Y-structure” pattern [Man06]. We here also give the generalization of the Y-structure pattern to simple SCMs.

**Proposition 9.9.1.** *Let  $M$  be a simple SCM with observed endogenous variables  $O = \{1, 2, 3, 4\} \subseteq V$  and no exogenous input variables ( $J = \emptyset$ ). Suppose that it is  $\sigma$ -faithful (w.r.t.  $O$ ). The following conditional (in)dependencies in the observational distribution  $P_M(X_1, X_2, X_3, X_4)$*

$$\begin{array}{lll} X_1 \not\perp\!\!\!\perp X_4, & X_2 \not\perp\!\!\!\perp X_4, & X_1 \perp\!\!\!\perp X_2, \\ X_1 \perp\!\!\!\perp X_4 \mid X_3, & X_2 \perp\!\!\!\perp X_4 \mid X_3, & X_1 \not\perp\!\!\!\perp X_2 \mid X_3, \end{array}$$

*imply that the observable graph  $G_O(M)$  must be one of the nine DMGs in Figure 25. Hence,*

1.  $X_4$  is not a cause of  $X_3$  according to  $M$ ;
2.  $X_3$  is a direct cause of  $X_4$  w.r.t.  $\{1, 2, 3, 4\}$  according to  $M$ ;
3.  $X_3$  and  $X_4$  do not have a common cause according to  $M$ ;
4. the causal effect of  $X_3$  on  $X_4$  satisfies:

$$P_M(X_4 \mid \text{do}(X_3)) = P_M(X_4 \mid X_3) \quad P_M(X_3)\text{-a.s.} \quad (57)$$

*Proof.* By using the global Markov property and the faithfulness assumption, one can check that the only (cyclic or acyclic) graphs that are compatible with the observed conditional independences are the ones in Figure 25. The statements now follow.  $\square$

This example illustrates how conditional independence patterns in the observational distribution allow one to infer certain features of the underlying causal model. This principle is exploited more generally by constraint-based methods, and implicitly, by

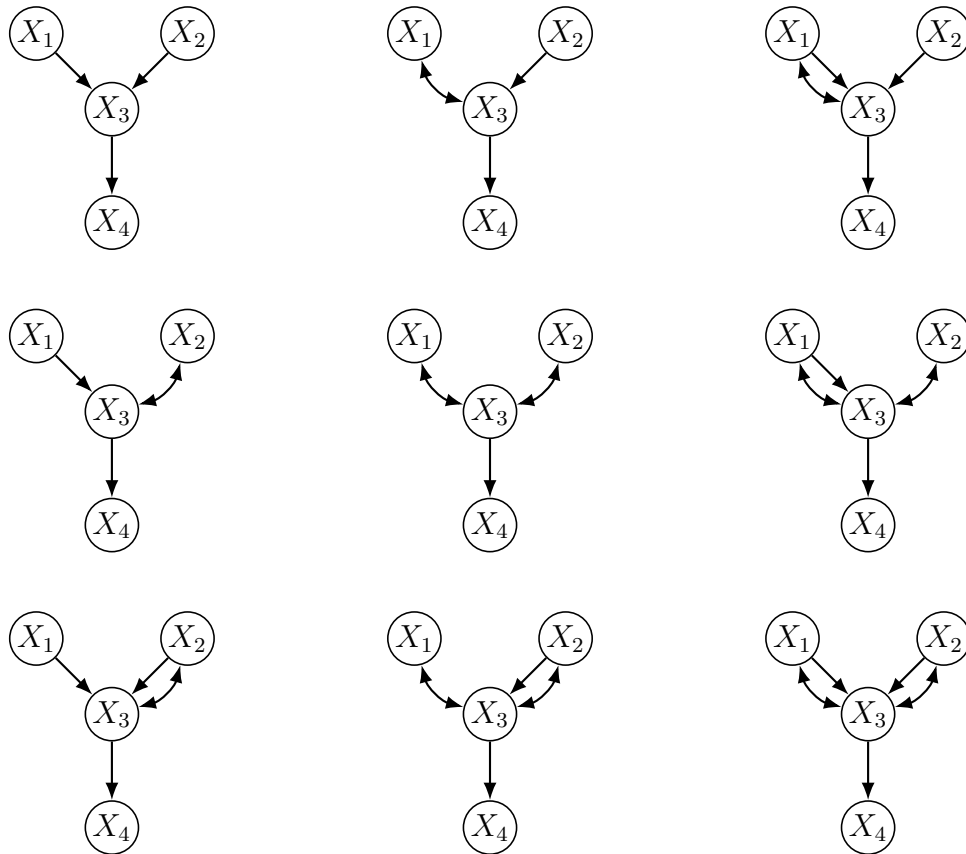


Figure 25: Observed causal graphs satisfying the “Y-structure” pattern on four variables.

score-based methods that optimize a penalized likelihood over (equivalence classes of) causal graphs. In Chapter 11 we will describe in detail one of the most sophisticated constraint-based causal discovery methods, Fast Causal Inference (FCI).

Typically, the graph cannot be completely identified from purely observational data. For example, in the Y-structure case, the conditional independences in the observational data do not allow to conclude whether the dependence between  $X_1$  and  $X_3$  is explained by  $X_1$  being a cause of  $X_3$ , or by  $X_1$  and  $X_3$  having a latent common cause, or both. However, under an appropriate faithfulness assumption, one can deduce the Markov equivalence class of the graph from the conditional independences in the observational data, i.e., the class of all CDMGs that induce the same separations.

Another disadvantage of causal discovery methods from purely observational data is that they typically need *very large* sample sizes and *strong assumptions* in order to work reliably (even on simulations).

## 9.10. Minimal Separating Sets, Minimal Connecting Sets

Minimal separating sets and minimal connecting sets are useful in that they give a relationship between certain separation properties of the graph and ancestral relations

in the graph [SMR99, CH11]. This can also be seen as a simple form of causal discovery.

**Definition 9.10.1.** Let  $X, Y, Z, S$  be sets of nodes in a CDMG  $G$  with input nodes  $J$  and output nodes  $V$ . We say that the minimal  $\sigma$ -separation

$$X \perp_G^\sigma Y \mid S \cup [Z]$$

holds if and only if

$$X \perp_G^\sigma Y \mid S \cup Z \quad \wedge \quad \forall Q \subsetneq Z : X \not\perp_G^\sigma Y \mid S \cup Q.$$

In words: all nodes in  $Z$  are required (in the context of the nodes in  $S$ ) to  $\sigma$ -separate  $X$  from  $Y$ . The minimal  $d$ -separation  $X \perp_G^d Y \mid S \cup [Z]$  is defined analogously.

Minimal separating sets imply the presence of certain ancestral relations (this generalizes a result of [SMR99]). But first we prove a little lemma.

**Lemma 9.10.2.** Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V$ . Let  $i, j \in J \cup V$  and  $Z \subseteq J \cup V$ . If  $\pi$  is a  $Z$ - $\sigma$ -open or  $Z$ - $d$ -open walk between  $i$  and  $j$  in  $G$ , then every node on  $\pi$  is in  $\text{Anc}_G(\{i, j\} \setminus J) \cup Z$ .

*Proof.* Suppose  $k$  is a node on  $\pi$ . Then either  $k$  is a collider, or there is a directed subwalk from  $k$  to a collider on  $\pi$ , or to an endnode of  $\pi$  that is not in  $J$ . In all cases,  $k \in \text{Anc}_G(\{i, j\} \setminus J) \cup Z$ . This holds for both  $d$ -separation and  $\sigma$ -separation.  $\square$

**Proposition 9.10.3.** Let  $\{x\}, \{y\}, S, Z$  be mutually disjoint sets of nodes in a CDMG  $G$  with input nodes  $J$  and output nodes  $V$ . Then:

$$x \perp_G^\sigma y \mid S \cup [Z] \implies Z \subseteq \text{Anc}_G(\{x, y\} \cup S).$$

A similar statement holds for  $d$ -separation.

*Proof.* Write  $A := \text{Anc}_G(\{x, y\} \cup S)$ . Let  $z \in Z$ . Suppose that  $z \notin A$ . Let  $Q = A \cap Z$ . Then  $z \notin Q$ , and therefore  $Q \subseteq Z \setminus \{z\}$ . Then there is a  $(Q \cup S)$ - $\sigma$ -open path  $\pi$  between  $x$  and  $\{y\} \cup J$  in  $G$ . Then every node on  $\pi$  is in  $\text{Anc}_G(\{x, y\} \setminus J) \cup Q \cup S$  (Lemma 9.10.2). Therefore, every node on  $\pi$  is in  $A$ . Hence no node in  $Z \setminus Q$  is on  $\pi$ . Therefore, adding  $(Z \setminus Q)$  to  $(Q \cup S)$  cannot  $\sigma$ -block  $\pi$ . Hence  $x \not\perp_G^\sigma y \mid Z \cup S$ . Contradiction.  $\square$

Similarly, we define minimal connections.

**Definition 9.10.4.** Let  $X, Y, Z, S$  be sets of nodes in a CDMG  $G$  with input nodes  $J$  and output nodes  $V$ . We say that the minimal  $\sigma$ -connection

$$X \not\perp_G^\sigma Y \mid S \cup [Z]$$

holds if and only if

$$X \not\perp_G^\sigma Y \mid S \cup Z \quad \wedge \quad \forall Q \subsetneq Z : X \perp_G^\sigma Y \mid S \cup Q.$$

In words: all nodes in  $Z$  are required (in the context of the nodes in  $S$ ) to  $\sigma$ -connect  $X$  with  $Y$ . The minimal  $d$ -connection  $X \not\perp_G^d Y \mid S \cup [Z]$  is defined analogously.

Note that despite the notation, a minimal connection is *not* the logical negation of a minimal separation.

Minimal connections imply the absence of certain ancestral relations:

**Proposition 9.10.5.** *Let  $\{x\}, \{y\}, S, \{z\}$  be mutually disjoint sets of nodes in a CDMG  $G$  with input nodes  $J$  and output nodes  $V$ . Then*

$$x \not\perp_G^\sigma y \mid S \cup \{z\} \implies z \notin \text{Anc}_G(\{x, y\} \cup S)$$

and a similar statement holds for  $d$ -separation.

*Proof.* There exists a  $S \cup \{z\}$ - $\sigma$ -open path between  $x$  and  $y \cup J$  in  $G$  that contains a collider in  $\text{Anc}_G(\{z\})$  that is not in  $\text{Anc}_G(S)$ . If  $z \in \text{Anc}_G(S)$  this would be a contradiction. If  $z \in \text{Anc}_G(x)$ , then we can consider the walk between  $x$  and  $y$  obtained from composing the subpath of the original path between  $y$  and the first collider (starting from  $y$ ) in  $\text{Anc}_G(\{z\}) \setminus \text{Anc}_G(S)$  with a directed path to  $z$  and then on to  $x$ , without passing through nodes in  $S$ . This walk between  $x$  and  $y$  must be  $\sigma$ -open given  $S$ , a contradiction. Similarly we obtain a contradiction if  $z \in \text{Anc}_G(y)$ .

The same proof works also for  $d$ -separation. □

## 10. Independence Testing

In this lecture, we will consider the following questions. How can we test whether...

- ... two random variables are independent?
- ... two random variables are conditionally independent given a third random variable?
- ... a random variable is independent of a non-random variable?
- ... a transitional random variable is conditionally independent of a transitional random variable, given another transitional random variable?

We will consider these questions only for the special case of finite categorical variables, i.e., variables that take values in finite spaces. In particular, we will discuss a test known as the  $G$  test. This has been defined in the literature for random variables, but we will extend it here to a general case involving transitional random variables (with “purely” random and “purely” non-random variables as special cases). We will state conditions under which the tests are asymptotically valid and consistent.

### 10.1. Marginal Independence for Categorical Random Variables

Consider two categorical random variables  $X, Y$  taking values in finite spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, with  $2 \leq |\mathcal{X}| < \infty$  and  $2 \leq |\mathcal{Y}| < \infty$ , and joint distribution  $P(X, Y)$ . We can represent the density in a table (assuming  $\mathcal{X} = \{1, \dots, k\}$  and  $\mathcal{Y} = \{1, \dots, l\}$ ):

	$Y = 1$	$Y = 2$	...	$Y = l$	
$X = 1$	$\theta_{11}$	$\theta_{12}$	...	$\theta_{1l}$	$\theta_{1+}$
$X = 2$	$\theta_{21}$	$\theta_{22}$	...	$\theta_{2l}$	$\theta_{2+}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$X = k$	$\theta_{k1}$	$\theta_{k2}$	...	$\theta_{kl}$	$\theta_{k+}$
	$\theta_{+1}$	$\theta_{+2}$	...	$\theta_{+l}$	$\theta_{++} = 1$

where we introduced the parameter  $\theta \in \Theta$  by setting  $\theta_{xy} = P(X = x, Y = y)$  for  $x \in \mathcal{X}, y \in \mathcal{Y}$ . We introduce here the convention that a “+” index denotes summation over that index, i.e.,

$$\theta_{+y} := \sum_{x \in \mathcal{X}} \theta_{xy}, \quad \theta_{x+} := \sum_{y \in \mathcal{Y}} \theta_{xy}, \quad \theta_{++} := \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \theta_{xy}.$$

For the parameter space we take the  $(|\mathcal{X}||\mathcal{Y}| - 1)$ -dimensional simplex:

$$\Theta := \left\{ \theta \in \prod_{(x,y) \in \mathcal{X} \times \mathcal{Y}} [0, 1] : \theta_{++} = 1 \right\}.$$

With Remark 2.5.23, we get:

$$X \underset{P(X,Y)}{\perp\!\!\!\perp} Y \iff P(X,Y) = P(X) \otimes P(Y),$$

where  $P(X)$  and  $P(Y)$  are the marginal distributions of  $P(X,Y)$ . In the discrete case we consider here, this holds if and only if

$$\forall x \in \mathcal{X}, y \in \mathcal{Y} : \theta_{xy} = \theta_{x+} \theta_{+y}.$$

The parameters satisfying this constraint form the allowed parameters under the null hypothesis of independence  $H_0 : X \perp\!\!\!\perp Y$ . We introduce the corresponding restricted parameter space

$$\Theta_0 := \{\theta \in \Theta : \theta_{xy} = \theta_{x+} \theta_{+y} \forall x \in \mathcal{X}, y \in \mathcal{Y}\} \subseteq \Theta.$$

We can also write the null hypothesis as  $H_0 : \theta \in \Theta_0$ . As alternative hypothesis we take that of dependence, i.e.,  $H_1 : X \not\perp\!\!\!\perp Y$ , or equivalently,  $H_1 : \theta \in \Theta_1$  with  $\Theta_1 := \Theta \setminus \Theta_0$ .

Suppose now that we have independent and identically distributed data  $(X_n, Y_n)_{n=1}^N$  with  $(X_n, Y_n) \sim P(X, Y | \theta)$  for all  $n = 1, \dots, N$ , with the “true” parameter  $\theta$  unknown. In other words, we assume for the joint distribution on the observed data

$$P((X_n, Y_n)_{n=1}^N | \theta) = \bigotimes_{n=1}^N P(X_n, Y_n | \theta),$$

where each  $P(X_n, Y_n | \theta)$  is a copy of the Markov kernel  $P(X, Y | \theta)$ . We define the *counts* as the number of observations with a given value  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ :

$$N_{xy} := \sum_{n=1}^N \mathbb{1}_{(x,y)}(X_n, Y_n).$$

We can represent them in a contingency table:

	$Y = 1$	$Y = 2$	$\dots$	$Y = l$	
$X = 1$	$N_{11}$	$N_{12}$	$\dots$	$N_{1l}$	$N_{1+}$
$X = 2$	$N_{21}$	$N_{22}$	$\dots$	$N_{2l}$	$N_{2+}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$X = k$	$N_{k1}$	$N_{k2}$	$\dots$	$N_{kl}$	$N_{k+}$
	$N_{+1}$	$N_{+2}$	$\dots$	$N_{+l}$	$N_{++} = N$

where we used a similar summation convention for the counts as for the parameters.

The classical frequentist procedure for deciding between two hypotheses  $H_0 : \theta \in \Theta_0$  and  $H_1 : \theta \in \Theta \setminus \Theta_0$  is as follows. One comes up with a *test statistic*  $T(D)$ , which is a function of the data  $D \sim P(D | \theta)$ , whose value should help us distinguish between the two hypotheses. We will consider here one-sided tests, where a large value of  $T(D)$  is in favor of  $H_1$  while a small value of  $T(D)$  is in favor of  $H_0$ . Then one chooses a

particular significance level  $\alpha \in (0, 1)$ . From the observed data  $d$ , one then calculates a corresponding  $p$ -value  $p(d)$ , which is the probability under the null hypothesis that the test statistic has the observed or a more extreme value. For the one-sided tests we will consider here, the  $p$ -value can be defined as

$$p(d) := \sup_{\theta \in \Theta_0} P(T(D) \geq T(d) \mid \theta).$$

Then, a decision is taken: if  $p(d) \leq \alpha$ , one considers this as sufficient evidence to reject  $H_0$  (and accept  $H_1$ ), while if  $p(d) > \alpha$ , one does not reject  $H_0$  as the evidence in the data is considered insufficient to do so. Often, the main desideratum is to control the probability of a Type I error (i.e., the error of incorrectly rejecting the null hypothesis), which can be achieved by choosing  $\alpha$  sufficiently small. Indeed, from the definition of the  $p$ -value it follows that:

$$\forall \alpha \in (0, 1) \forall \theta \in \Theta_0 : P(p(D) \leq \alpha \mid \theta) \leq \alpha.$$

For causal discovery, however, we need a more symmetric treatment of the two hypotheses, as there we require both the probability of a Type I error and of a Type II error (i.e., the error of incorrectly rejecting the alternative hypothesis) to be small. Before we investigate this tradeoff, let us first propose a concrete test statistic for the case at hand and obtain an approximate expression for the corresponding  $p$ -value.

Here we will work out the details of the *likelihood ratio test*, which for this particular case is also known as the *G test*. We start by writing down the likelihood of the data:

$$P((X_n, Y_n)_{n=1}^N \mid \theta) = \prod_{n=1}^N \theta_{X_n Y_n} = \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \theta_{xy}^{N_{xy}}$$

where we used the counts as a sufficient statistic of the data. This is a multinomial distribution with parameters  $(\theta_{xy})_{x \in \mathcal{X}, y \in \mathcal{Y}}$  and  $N$ . Maximizing the likelihood with respect to the parameters  $\theta \in \Theta$ , we obtain the well-known maximum likelihood estimator

$$\hat{\theta}_{xy} = \frac{N_{xy}}{N},$$

i.e., the fractions of the different outcomes in the data. Under the null hypothesis  $H_0$ ,  $\theta_{xy} = \theta_{x+} \theta_{+y}$ , and the likelihood factorizes:

$$\begin{aligned} \theta \in \Theta_0 \implies P((X_n, Y_n)_{n=1}^N \mid \theta) &= \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} (\theta_{x+} \theta_{+y})^{N_{xy}} \\ &= \left( \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \theta_{x+}^{N_{xy}} \right) \left( \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \theta_{+y}^{N_{xy}} \right) \\ &= \left( \prod_{x \in \mathcal{X}} \theta_{x+}^{N_{x+}} \right) \left( \prod_{y \in \mathcal{Y}} \theta_{+y}^{N_{+y}} \right). \end{aligned}$$

This is just the product of two independent multinomial distributions with (variationally independent) parameters  $(\theta_{x+})_{x \in \mathcal{X}}$  and  $(\theta_{+y})_{y \in \mathcal{Y}}$  (and  $N$ ), respectively. Hence, the restricted maximum likelihood estimator under  $H_0$  is

$$\hat{\theta}_{xy}^0 = \hat{\theta}_{x+}^0 \hat{\theta}_{+y}^0 = \frac{N_{x+}}{N} \frac{N_{+y}}{N}.$$

The likelihood ratio is obtained by dividing the likelihood for  $\hat{\theta}$  by the likelihood for  $\hat{\theta}^0$ :

$$\begin{aligned} \frac{\sup_{\theta \in \Theta_0 \cup \Theta_1} P((X_n, Y_n)_{n=1}^N | \theta)}{\sup_{\theta \in \Theta_0} P((X_n, Y_n)_{n=1}^N | \theta)} &= \frac{P((X_n, Y_n)_{n=1}^N | \hat{\theta})}{P((X_n, Y_n)_{n=1}^N | \hat{\theta}^0)} = \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \left( \frac{\hat{\theta}_{xy}}{\hat{\theta}_{x+}^0 \hat{\theta}_{+y}^0} \right)^{N_{xy}} \\ &= \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \left( \frac{N_{xy} N}{N_{x+} N_{+y}} \right)^{N_{xy}}. \end{aligned} \quad (58)$$

The likelihood ratio test statistic is defined as 2 times the natural logarithm of this ratio:

$$G_N := 2 \log \frac{\sup_{\theta \in \Theta_0 \cup \Theta_1} P((X_n, Y_n)_{n=1}^N | \theta)}{\sup_{\theta \in \Theta_0} P((X_n, Y_n)_{n=1}^N | \theta)} = 2 \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} N_{xy} \log \frac{N_{xy} N}{N_{x+} N_{+y}}. \quad (59)$$

Since counts can be zero, one should interpret  $0 \log \frac{0}{n}$  in this expression as 0 (for  $n \in \mathbb{N}$ ).

We will now consider the asymptotic behavior of the test statistic under the null hypothesis. This will yield an approximation for the  $p$ -value that we can use also for finite samples. As a simplifying assumption, we will henceforth assume that all probabilities are positive,<sup>49</sup> i.e.,

$$\theta_{xy} > 0 \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}. \quad (60)$$

**Proposition 10.1.1.** *Under  $H_0 : X \perp\!\!\!\perp Y$ , and with regularity assumption (60),*

$$G_N \rightsquigarrow \chi_\nu^2$$

with  $\nu = (|\mathcal{X}| - 1)(|\mathcal{Y}| - 1)$  as sample size  $N \rightarrow \infty$ . In words, the test statistic  $G_N$  converges in distribution<sup>50</sup> as  $N \rightarrow \infty$  to a chi-squared distribution with  $\nu$  degrees of freedom.<sup>51</sup>

*Proof.* This is a direct application of Theorem 4.43 in [BJvdV17], for which a proof is provided in Chapter 16 in [vdV98]. One has to be careful here to use a different parameterization—in terms of (variationally) independent parameters—i.e., such

<sup>49</sup>The singularities for vanishing values of  $\theta_{xy}$  can be dealt with, but require special attention. For simplicity we study only the regular case here.

<sup>50</sup>We say that a sequence of real-valued random variables  $X_1, X_2, \dots$  converges in distribution to  $X_\infty$ , and write  $X_n \rightsquigarrow X_\infty$ , if  $P(X_n \leq x) \rightarrow P(X_\infty \leq x)$  for all  $x \in \mathbb{R}$  such that  $\xi \mapsto P(X_\infty \leq \xi)$  is continuous at  $x$ .

<sup>51</sup>The chi-square distribution with  $\nu$  degrees of freedom is defined as the distribution of a sum of squares of  $\nu$  independent standard normal random variables, i.e., of  $\sum_{i=1}^\nu Z_i^2$  where  $Z_i \sim N(0, 1)$  are i.i.d..



that the parameter space contains an open part of  $\mathbb{R}^{|\mathcal{X}||\mathcal{Y}|-1}$ , when calculating the score function and the Fisher information matrix when checking the regularity conditions. For example, one can choose a pair  $(k, l) \in \mathcal{X} \times \mathcal{Y}$  and take parameters  $\theta_{x,y} = \vartheta_{x,y}$  for  $x \neq k$  or  $y \neq l$ , and  $\theta_{k,l} = 1 - \sum_{(x,y) \neq (k,l)} \vartheta_{x,y}$ . The dimensionality of  $\Theta$  is  $|\mathcal{X}||\mathcal{Y}| - 1$ , while that of  $\Theta_0$  is  $(|\mathcal{X}| - 1) + (|\mathcal{Y}| - 1)$ . The degrees of freedom of the asymptotic chi-square distribution is the difference of the two, i.e.,  $\nu = (|\mathcal{X}||\mathcal{Y}| - 1) - ((|\mathcal{X}| - 1) + (|\mathcal{Y}| - 1)) = (|\mathcal{X}| - 1)(|\mathcal{Y}| - 1)$ .

An alternative proof will be provided later in a more general setting (see Proposition 10.3.2).  $\square$

One therefore obtains an approximate level  $\alpha$  test (i.e., a test with Type I error asymptotically upper bounded by  $\alpha$ ) by rejecting  $H_0$  when  $G_N \geq \chi_{\nu, 1-\alpha}^2$ . Here,  $\chi_{\nu, 1-\alpha}^2 := F_{\chi_\nu^2}^{-1}(1 - \alpha)$  is the upper  $\alpha$  quantile of the  $\chi^2$ -distribution with  $\nu$  degrees of freedom, with  $F_{\chi_\nu^2}$  the corresponding distribution function (cumulative density function) and  $F_{\chi_\nu^2}^{-1}$  its inverse (i.e., the quantile function). Indeed, if  $\theta \in \Theta_0$ , then  $P(G_N \geq \chi_{\nu, 1-\alpha}^2) \rightarrow \alpha$ , for any  $\alpha \in (0, 1)$ . Since

$$G_N \geq \chi_{\nu, 1-\alpha}^2 \iff G_N \geq F_{\chi_\nu^2}^{-1}(1 - \alpha) \iff F_{\chi_\nu^2}(G_N) \geq 1 - \alpha \iff 1 - F_{\chi_\nu^2}(G_N) \leq \alpha,$$

the corresponding approximate  $p$ -value is  $1 - F_{\chi_\nu^2}(G_N)$ ; if this is smaller than or equal to the chosen threshold  $\alpha$ , we reject  $H_0$ . This test is called the  $G$ -test.

But what about the Type II error? If we let the sample size  $N$  grow, we would hope that the probability of a wrong test result becomes arbitrarily small, and vanishes in the limit  $N \rightarrow \infty$ .

**Definition 10.1.2.** *A (conditional) independence test is called consistent if the probabilities of both Type I and Type II errors converge to 0, no matter what the true parameter value is.*

To obtain consistency, it is not an option to just control Type I error at a fixed level  $\alpha$ ; instead, one has to use a level  $\alpha_N$  that depends on the sample size  $N$ , and converges to 0 (implying that Type I error converges to 0). However, because of the tradeoff between Type I and Type II errors, the rate at which  $\alpha_N$  converges to 0 has to be chosen carefully in order to be able to guarantee that also Type II error vanishes asymptotically. As we shall see, the convergence rate of  $\alpha_N$  should be chosen sufficiently slow.

While it is often easier to calculate the Type I error than the Type II error of a test, in this case we can actually analyze the asymptotic behavior of the test statistic under the alternative hypothesis  $H_1$ . Define

$$\hat{I}_N := \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \hat{\theta}_{xy} \log \frac{\hat{\theta}_{xy}}{\hat{\theta}_x \hat{\theta}_y} = \frac{G_N}{2N}$$

where we used that  $\hat{\theta}_{x+}^0 = \hat{\theta}_{x+}$  and  $\hat{\theta}_{+y}^0 = \hat{\theta}_{+y}$ . This is an estimator (the so-called “plug-in

estimator”  $I(\hat{\theta})$  of the mutual information  $I(X; Y)$ :

$$\begin{aligned} I(\theta) &:= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \theta_{xy} \log \frac{\theta_{xy}}{\theta_{x+} \theta_{y+}} \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P(X = x, Y = y) \log \frac{P(X = x, Y = y)}{P(X = x)P(Y = y)} =: I(X; Y). \end{aligned}$$

With Jensen’s inequality, one can show that  $I(X; Y) \geq 0$ , and that  $I(X; Y) = 0 \iff X \perp\!\!\!\perp Y$ . Note further that the function  $\Theta \rightarrow [0, \infty) : \theta \mapsto I(\theta)$  is continuous.

With this observation, we can prove the asymptotic consistency of the  $G$ -test under assumptions on the critical values used for deciding between  $H_0$  and  $H_1$ .

**Corollary 10.1.3.** *Consider an infinite sequence of  $G$  tests performed on the first  $N$  samples of an infinitely large data set  $(X_n, Y_n)_{n=1}^\infty$ , where one accepts  $H_1 : X \not\perp\!\!\!\perp Y$  if  $G_N \geq \tau_N$ , and otherwise accepts  $H_0 : X \perp\!\!\!\perp Y$ , for some given sequence of thresholds  $\tau_N$ . Under the regularity assumption (60), this sequence of tests is asymptotically consistent if  $\tau_N \rightarrow \infty$  but  $\tau_N/N \rightarrow 0$ .*

*Proof.* We start by a simple application of the strong law of large numbers. Let  $\theta \in \Theta$ . Since the  $(X_n, Y_n)$  are assumed to be i.i.d., and

$$\mathbb{E}(\mathbb{1}_{(x,y)}(X_n, Y_n)) = \theta_{xy}$$

for all  $x \in \mathcal{X}, y \in \mathcal{Y}$ , we conclude that  $N_{xy}/N \xrightarrow{\text{a.s.}} \theta_{xy}$  for all  $x \in \mathcal{X}, y \in \mathcal{Y}$  by the strong law of large numbers.<sup>52</sup> Hence  $\hat{\theta}_{xy} \xrightarrow{\text{a.s.}} \theta_{xy}$ . Hence, also  $\hat{\theta}_{+x} \xrightarrow{\text{a.s.}} \theta_{+x}$  and  $\hat{\theta}_{y+} \xrightarrow{\text{a.s.}} \theta_{y+}$ . Furthermore, by continuity,  $I(\hat{\theta}) \xrightarrow{\text{a.s.}} I(\theta)$ . Hence,  $G_N/N \xrightarrow{\text{a.s.}} 2I(\theta)$ .

Under  $H_1$ , we have  $I(\theta) > 0$ , and since by assumption  $\tau_N/N \rightarrow 0$ ,  $\mathbb{1}_{G_N < \tau_N} \xrightarrow{\text{a.s.}} 0$ . Since a.s. convergence implies convergence in probability,

$$\theta \in \Theta_1 \implies P(G_N < \tau_N) \rightarrow 0.$$

Thus, the probability of a Type II error vanishes asymptotically.

This same approach doesn’t work for the Type I error. The reason is that even though we assume  $\tau_N \rightarrow \infty$ , and we know that  $G_N/N \xrightarrow{\text{a.s.}} 0$  under  $H_0$ , this does not suffice to conclude anything about the probability of the event  $G_N \geq \tau_N$ . But we can make use of Proposition 10.1.1, which states that  $G_N \rightsquigarrow \chi_\nu^2$  under  $H_0$ . Since the distribution function of  $\chi_\nu^2$  is continuous, this implies uniform convergence of the distribution functions:

$$\sup_{x \in \mathbb{R}} |F_{G_N}(x) - F_{\chi_\nu^2}(x)| \rightarrow 0.$$

Hence

$$|F_{G_N}(\tau_N) - F_{\chi_\nu^2}(\tau_N)| \leq \sup_{x \in \mathbb{R}} |F_{G_N}(x) - F_{\chi_\nu^2}(x)| \rightarrow 0.$$

<sup>52</sup>The convergence is “almost surely”, i.e.,  $N_{xy}/N \xrightarrow{\text{a.s.}} \theta_{xy}$  means that  $P(N_{xy}/N \rightarrow \theta_{xy}) = 1$ .

Since  $\tau_N \rightarrow \infty$ ,  $F_{\chi^2_\nu}(\tau_N) \rightarrow 1$ . Hence, also  $F_{G_N}(\tau_N) \rightarrow 1$ . We conclude that

$$\theta \in \Theta_0 \implies P(G_N \geq \tau_N) \rightarrow 0,$$

i.e., the probability of a Type I error converges to 0.  $\square$

While one traditionally focuses mostly on Type I error control, in causal discovery we are more interested in having both small Type I and Type II error. In order to achieve this (at least asymptotically, i.e., for sufficiently large sample sizes), we can thus make use of a sequence of thresholds that satisfies the assumptions in the corollary. In terms of  $p$ -values, this means that to bound the Type I error, a fixed critical value  $\alpha$  suffices, but for consistency we let  $\alpha_N \rightarrow 0$  with a rate such that  $\chi^2_{\nu, 1-\alpha_N}/N \rightarrow 0$ .

While for a finite sample, we can give guarantees (at least approximately) on the Type I error, it will often be impossible to provide guarantees on the Type II error without making strong assumptions on the parameters. Indeed, since the mutual information  $I(X; Y)$  (a measure of the dependence of  $X$  and  $Y$ ) can be arbitrarily close to zero for weakly dependent  $X$  and  $Y$ , one cannot know in advance how many samples will be needed to be able to distinguish it from an independence.<sup>53</sup>

## 10.2. Conditional Independence for Categorical Random Variables

We now extend the  $G$  test to a conditional independence test that we will refer to as the conditional  $G$  test.

Consider three categorical random variables  $X, Y, Z$  taking values in spaces  $\mathcal{X}, \mathcal{Y}$  and  $\mathcal{Z}$ , respectively (with  $2 \leq |\mathcal{X}| < \infty$ ,  $2 \leq |\mathcal{Y}| < \infty$  and  $1 \leq |\mathcal{Z}| < \infty$ ) and joint distribution  $P(X, Y, Z)$ . With Remark 2.5.23, we get (because finite spaces are standard):

$$\begin{aligned} X \underset{P(X, Y, Z)}{\perp\!\!\!\perp} Y \mid Z &\iff P(X, Y, Z) = P(X \mid Z) \otimes P(Y, Z) \\ &\iff P(X, Y \mid Z) = P(X \mid Z) \otimes P(Y \mid Z) \quad P(Z)\text{-a.s.} \\ &\iff \forall z \in \mathcal{Z} : [P(Z = z) > 0 \implies \\ &\quad P(X, Y \mid Z = z) = P(X \mid Z = z)P(Y \mid Z = z)] \\ &\iff \forall z \in \mathcal{Z} : [P(Z = z) > 0 \implies X \underset{P(X, Y \mid Z = z)}{\perp\!\!\!\perp} Y]. \end{aligned}$$

This suggests that we can make use of an independence test for two categorical variables on each “stratum” corresponding to conditioning on a specific value  $Z = z$  that has positive probability to occur.

We parameterize the conditional kernel  $P(X, Y \mid Z)$  in terms of parameters  $(\theta_{xy|z})_{x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}}$  which live in space

$$\Theta_{XY|Z} := \left\{ \theta \in \prod_{x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}} [0, 1] : \forall z \in \mathcal{Z} \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \theta_{xy|z} = 1 \right\}.$$

<sup>53</sup>This is referred to as the lack of “uniformly consistent” (conditional) independence tests.

With the summation convention, we can write the normalization condition as  $\theta_{++|z} = 1$  for all  $z \in \mathcal{Z}$ . For those  $z \in \mathcal{Z}$  with  $P(Z = z) > 0$ , we have

$$\frac{P(X = x, Y = y, Z = z)}{P(Z = z)} = P(X = x, Y = y \mid Z = z) = \theta_{xy|z}.$$

We also parameterize the marginal distribution  $P(Z)$  in terms of parameters  $(\theta_z)_{z \in \mathcal{Z}}$  which live in space

$$\Theta_Z := \{\theta \in \prod_{z \in \mathcal{Z}} [0, 1] : \sum_{z \in \mathcal{Z}} \theta_z = 1\}.$$

Any joint distribution of  $X, Y$  and  $Z$  can then be parameterized as

$$P(X = x, Y = y, Z = z \mid \theta) = \theta_z \theta_{xy|z},$$

with parameter space

$$\Theta := \Theta_Z \times \Theta_{XY|Z}.$$

We formulate the null hypothesis  $H_0 : X \perp\!\!\!\perp Y \mid Z$  of independence in terms of the parameters as

$$\forall x \in \mathcal{X}, y \in \mathcal{Y}, \forall z \in \mathcal{Z} : \theta_{xy|z} = \theta_{x+|z} \theta_{+y|z}$$

(for convenience, we have strengthened it a bit; strictly speaking, we only need this relation to hold for all  $z \in \mathcal{Z}$  with  $\theta_z > 0$ ; however, since the data will not convey any information on  $\theta_{xy|z}$  for such  $z$ , this does not matter). The corresponding restricted parameter space is

$$\Theta_{XY|Z}^0 := \{\theta \in \Theta : \theta_{xy|z} = \theta_{x+|z} \theta_{+y|z} \forall x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}\}.$$

We can then also write the null hypothesis as  $H_0 : \theta \in \Theta_0$  with  $\Theta_0 := \Theta_Z \times \Theta_{XY|Z}^0$ . As alternative hypothesis we take that of dependence, i.e.,  $H_1 : X \not\perp\!\!\!\perp Y \mid Z$ , or equivalently,  $H_1 : \theta \in \Theta_1$ , where  $\Theta_1 := \Theta_Z \times \Theta_{XY|Z}^1$  with  $\Theta_{XY|Z}^1 := \Theta_{XY|Z} \setminus \Theta_{XY|Z}^0$ .

Suppose now that we have independent and identically distributed data  $(X_n, Y_n, Z_n)_{n=1}^N$  with  $(X_n, Y_n, Z_n) \sim P(X, Y, Z \mid \theta)$  for all  $n = 1, \dots, N$ , with the “true” parameter  $\theta \in \Theta$  unknown. In other words, we assume for the joint distribution on the observed data

$$P((X_n, Y_n, Z_n)_{n=1}^N \mid \theta) = \bigotimes_{n=1}^N P(X_n, Y_n, Z_n \mid \theta),$$

where each  $P(X_n, Y_n, Z_n \mid \theta)$  is a copy of the Markov kernel  $P(X, Y, Z \mid \theta)$ . We define the *counts* as the number of observations with a given value  $(x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$ :

$$N_{xyz} := \sum_{n=1}^N \mathbb{1}_{(x,y,z)}(X_n, Y_n, Z_n).$$

We again work out the details of the likelihood ratio test, and start by writing down the likelihood of the data:

$$\begin{aligned} P((X_n, Y_n, Z_n)_{n=1}^N | \theta) &= \prod_{n=1}^N (\theta_{X_n, Y_n | Z_n} \theta_{Z_n}) = \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \prod_{z \in \mathcal{Z}} (\theta_{xy|z}^{N_{xyz}} \theta_z^{N_{xyz}}) \\ &= \left( \prod_{z \in \mathcal{Z}} \theta_z^{N_{++z}} \right) \left( \prod_{z \in \mathcal{Z}} \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \theta_{xy|z}^{N_{xyz}} \right) \end{aligned}$$

where we used the counts as a sufficient statistic of the data. We recognize the first factor as the likelihood of a multinomial distribution with parameters  $(\theta_z)_{z \in \mathcal{Z}}$  and  $N$ . The second factor is a product of the likelihoods of multinomial distributions with parameters  $(\theta_{xy|z})_{x \in \mathcal{X}, y \in \mathcal{Y}}$  and  $N_{++z}$ , for each  $z \in \mathcal{Z}$ . Maximizing the likelihood with respect to the parameters  $\theta \in \Theta$ , we obtain the maximum likelihood estimator

$$\left( \hat{\theta}_{xy|z}, \hat{\theta}_z \right) = \left( \frac{N_{xyz}}{N_{++z}}, \frac{N_{++z}}{N} \right).$$

Under the null hypothesis  $H_0$ ,  $\theta_{xy|z} = \theta_{x+|z} \theta_{+y|z}$ , and the likelihood factorizes over  $X$  and  $Y$ :

$$\begin{aligned} \theta \in \Theta_0 &\implies P((X_n, Y_n, Z_n)_{n=1}^N | \theta) = \left( \prod_{z \in \mathcal{Z}} \theta_z^{N_{++z}} \right) \prod_{z \in \mathcal{Z}} \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} (\theta_{x+|z} \theta_{+y|z})^{N_{xyz}} \\ &= \left( \prod_{z \in \mathcal{Z}} \theta_z^{N_{++z}} \right) \prod_{z \in \mathcal{Z}} \left( \prod_{x \in \mathcal{X}} \theta_{x+|z}^{N_{x+z}} \right) \left( \prod_{y \in \mathcal{Y}} \theta_{+y|z}^{N_{+yz}} \right). \end{aligned}$$

The restricted maximum likelihood estimator under  $H_0$  is

$$\left( \hat{\theta}_{xy|z}^0, \hat{\theta}_z^0 \right) = \left( \hat{\theta}_{x+|z}^0 \hat{\theta}_{+y|z}^0, \hat{\theta}_z^0 \right) = \left( \frac{N_{x+z} N_{+yz}}{N_{++z}^2}, \frac{N_{++z}}{N} \right).$$

The likelihood ratio is obtained by dividing the likelihood for  $\hat{\theta}$  by the likelihood for  $\hat{\theta}^0$ :

$$\begin{aligned} \frac{\sup_{\theta \in \Theta_0 \cup \Theta_1} P((X_n, Y_n, Z_n)_{n=1}^N | \theta)}{\sup_{\theta \in \Theta_0} P((X_n, Y_n, Z_n)_{n=1}^N | \theta)} &= \frac{P((X_n, Y_n, Z_n)_{n=1}^N | \hat{\theta})}{P((X_n, Y_n, Z_n)_{n=1}^N | \hat{\theta}^0)} = \prod_{z \in \mathcal{Z}} \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \left( \frac{\hat{\theta}_{xy|z}}{\hat{\theta}_{x+|z}^0 \hat{\theta}_{+y|z}^0} \right)^{N_{xyz}} \\ &= \prod_{z \in \mathcal{Z}} \prod_{x \in \mathcal{X}} \prod_{y \in \mathcal{Y}} \left( \frac{N_{xyz} N_{++z}}{N_{x+z} N_{+yz}} \right)^{N_{xyz}}, \end{aligned}$$

where the factors involving the marginal  $P(Z)$  cancel out. The likelihood ratio test statistic is defined as 2 times the natural logarithm of this ratio:

$$G_N := 2 \log \frac{\sup_{\theta \in \Theta_0 \cup \Theta_1} P((X_n, Y_n, Z_n)_{n=1}^N | \theta)}{\sup_{\theta \in \Theta_0} P((X_n, Y_n, Z_n)_{n=1}^N | \theta)} = 2 \sum_{z \in \mathcal{Z}} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} N_{xyz} \log \frac{N_{xyz} N_{++z}}{N_{x+z} N_{+yz}}. \quad (61)$$

We will now consider the asymptotic behavior of the test statistic under both hypotheses. As a simplifying assumption, we will henceforth assume that all probabilities are positive, i.e.,

$$\begin{cases} \theta_z > 0 & \forall z \in \mathcal{Z}, \text{ and} \\ \theta_{xy|z} > 0 & \forall x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}. \end{cases} \quad (62)$$

**Proposition 10.2.1.** *Under  $H_0 : X \perp\!\!\!\perp Y | Z$ , and with regularity assumption (62)*

$$G_N \rightsquigarrow \chi_\nu^2$$

with  $\nu = |\mathcal{Z}|(|\mathcal{X}| - 1)(|\mathcal{Y}| - 1)$  as sample size  $N \rightarrow \infty$ . In words, the test statistic  $G_N$  converges in distribution as  $N \rightarrow \infty$  to a chi-squared distribution with  $\nu$  degrees of freedom.

*Proof.* This is analogous to the proof of Proposition 10.1.1. The dimensionality of  $\Theta$  is  $|\mathcal{Z}|(|\mathcal{X}||\mathcal{Y}| - 1) + (|\mathcal{Z}| - 1)$ , while that of  $\Theta_0$  is  $|\mathcal{Z}|((|\mathcal{X}| - 1) + (|\mathcal{Y}| - 1)) + (|\mathcal{Z}| - 1)$ . The degrees of freedom of the asymptotic chi-square distribution is the difference of the two, i.e.,  $\nu = |\mathcal{Z}|(|\mathcal{X}| - 1)(|\mathcal{Y}| - 1)$ .  $\square$

One therefore obtains an approximate level  $\alpha$  test (i.e., a test with Type I error asymptotically upper bounded by  $\alpha$ ) by rejecting  $H_0$  when  $G_N \geq \chi_{\nu, 1-\alpha}^2$ .

Define

$$\hat{I}_N := \sum_{z \in \mathcal{Z}} \hat{\theta}_z \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \hat{\theta}_{xy|z} \log \frac{\hat{\theta}_{xy|z}}{\hat{\theta}_{x+|z} \hat{\theta}_{+y|z}} = \frac{G_N}{2N}$$

where we used that  $\hat{\theta}_{x+|z}^0 = \hat{\theta}_{x+|z}$  and  $\hat{\theta}_{+y|z}^0 = \hat{\theta}_{+y|z}$ . This is a plug-in estimator of the conditional mutual information  $I(X; Y|Z)$ :

$$\begin{aligned} I(\theta) &:= \sum_{z \in \mathcal{Z}} \theta_z \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \theta_{xy|z} \log \frac{\theta_{xy|z}}{\theta_{x+|z} \theta_{+y|z}} \\ &= \sum_{z \in \mathcal{Z}} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P(X = x, Y = y, Z = z) \log \frac{P(X = x, Y = y | Z = z)}{P(X = x | Z = z) P(Y = y | Z = z)} \\ &=: I(X; Y|Z). \end{aligned}$$

With Jensen's inequality, one can show that  $I(X; Y|Z) \geq 0$ , and that  $I(X; Y|Z) = 0 \iff X \perp\!\!\!\perp Y | Z$ . Note further that the function  $\Theta \rightarrow [0, \infty) : \theta \mapsto I(\theta)$  is continuous.

With this observation, we can prove the asymptotic consistency of the conditional  $G$ -test under assumptions on the critical values used for deciding between  $H_0$  and  $H_1$ .

**Corollary 10.2.2.** *Consider an infinite sequence of conditional  $G$  tests performed on the first  $N$  samples of an infinitely large data set  $(X_n, Y_n, Z_n)_{n=1}^\infty$ , where one accepts  $H_1 : X \not\perp\!\!\!\perp Y | Z$  if  $G_N \geq \tau_N$ , and otherwise accepts  $H_0 : X \perp\!\!\!\perp Y | Z$ , for some given sequence of thresholds  $\tau_N$ . Under the regularity assumption (62), this sequence is asymptotically consistent if  $\tau_N \rightarrow \infty$  but  $\tau_N/N \rightarrow 0$ .*

*Proof.* This is very similar to the proof of Corollary 10.1.3.

We again apply the strong law of large numbers. Let  $\theta \in \Theta$ . Since  $(X_n, Y_n, Z_n)$  are assumed to be i.i.d., and

$$\mathbb{E}(\mathbb{1}_{(x,y,z)}(X_n, Y_n, Z_n)) = \theta_{xy|z}\theta_z$$

for all  $x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}$ , we conclude that  $N_{xyz}/N \xrightarrow{a.s.} \theta_{xy|z}\theta_z$  for all  $x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}$  by the strong law of large numbers. Hence, also  $N_{++z}/N \xrightarrow{a.s.} \theta_z$  for all  $z \in \mathcal{Z}$ . Hence, using (62),  $\hat{\theta}_z \xrightarrow{a.s.} \theta_z$ ,  $\hat{\theta}_{xy|z} \xrightarrow{a.s.} \theta_{xy|z}$ ,  $\hat{\theta}_{+x|z} \xrightarrow{a.s.} \theta_{+x|z}$ , and  $\hat{\theta}_{y+|z} \xrightarrow{a.s.} \theta_{y+|z}$ . By continuity,  $I(\hat{\theta}) \xrightarrow{a.s.} I(\theta)$ . Hence,  $G_N/N \xrightarrow{a.s.} 2I(\theta)$ .

We can now reason analogously as in the proof of Corollary 10.1.3 to conclude that the probability of a Type II error vanishes asymptotically.

For an asymptotic estimate of the probability of a Type I error, we can make use of Proposition 10.2.1, which states that  $G_N \rightsquigarrow \chi_\nu^2$ . This part of the proof is identical to the corresponding part of the proof of Corollary 10.1.3.  $\square$

### 10.3. Marginal Independence of a Random and a Non-Random Variable

Consider two variables  $X, C$  taking values in finite spaces  $\mathcal{X}, \mathcal{C}$ , respectively, (i.e., with  $2 \leq |\mathcal{X}| < \infty$  and  $2 \leq |\mathcal{C}| < \infty$ ). Assume that  $X$  is a random variable, while  $C$  is an input variable, and consider a Markov kernel  $K(X | C) : \mathcal{C} \dashrightarrow \mathcal{X}$ . We will derive a test for the independence

$$X \underset{K(X|C)}{\perp\!\!\!\perp} C.$$

With Definition 2.5.17, this independence holds if and only if there exists a Markov kernel  $Q(X) : * \dashrightarrow \mathcal{X}$  such that:

$$K(X | C) = Q(X).$$

The latter means that

$$K(X | C = c) = Q(X) \tag{63}$$

for all  $c \in \mathcal{C}$ .

Suppose we obtain data  $(X_n, c_n)_{n=1}^N$  such that the  $X_n$  are conditionally independent and identically distributed given  $c_n$  for all  $n = 1, \dots, N$ . In other words, we assume the data is sampled from the following Markov kernel:

$$K((X_n)_{n=1}^N | (c_n)_{n=1}^N) = \bigotimes_{n=1}^N K(X_n | C_n = c_n),$$

where each  $K(X_n | C_n)$  is a copy of the (“true” but unknown kernel)  $K(X | C)$ .

Note that this is a weaker assumption regarding the sampling scheme than we would have made if  $C$  were random. In particular, we make no assumption at all regarding how the sequence of values  $c_1, c_2, \dots, c_N$  is chosen. It could be a sequence like

0, 1, 0, 1, 0, 1, 0, 1,  $\dots$ , for example, which would be (if sufficiently long) very unlikely to occur if all  $c_n$  would be independently sampled from some distribution. This extends the possible experimental designs that we can handle to include for example randomized controlled trials in which the protocol is such that a certain prespecified number  $N_{+|0}$  of subjects enters the control group, and a certain prespecified number  $N_{+|1}$  enters the treatment group. If the values of  $c_n$  were chosen i.i.d. with a coin flip, then it would be very unlikely that this assignment satisfies the protocol. Considering  $C$  to be an exogenous input variable instead (with values that are not necessarily randomly assigned), allows us to test for the independence of outcome  $X$  and treatment  $C$  under a broader range of experimental protocols.

We define the *counts* as the number of observations with a given value  $(x, c) \in \mathcal{X} \times \mathcal{C}$ :

$$N_{x|c} := \sum_{n=1}^N \mathbb{1}_{(x,c)}(X_n, c_n).$$

We will take  $H_0 : X \perp\!\!\!\perp C$  as the null hypothesis of a frequentist test for the independence of  $X$  and  $C$ . We parameterize the Markov kernel  $K(X | C)$  in terms of parameters  $(\theta_{x|c})_{x \in \mathcal{X}, c \in \mathcal{C}} := K(X = x | C = c)$  in a space

$$\Theta := \left\{ \theta \in \prod_{x \in \mathcal{X}, c \in \mathcal{C}} [0, 1] : \forall c \in \mathcal{C} \sum_{x \in \mathcal{X}} \theta_{x|c} = 1 \right\}.$$

With the summation convention, we can write the normalization condition as  $\theta_{+|c} = 1$  for all  $c \in \mathcal{C}$ . The null hypothesis  $H_0 : X \perp\!\!\!\perp C$ , equivalent to (63), can be expressed in terms of the parameters as  $H_0 : \theta \in \Theta_0$ , where we introduced the restricted parameter space

$$\Theta_0 := \left\{ \theta \in \Theta : \forall x \in \mathcal{X}, c \in \mathcal{C}, c' \in \mathcal{C} : \theta_{x|c} = \theta_{x|c'} \right\}.$$

As alternative hypothesis we will take  $H_1 : X \not\perp\!\!\!\perp C$ , the negation of the null hypothesis, i.e.,  $H_1 : \theta \in \Theta_1$  with  $\Theta_1 := \Theta \setminus \Theta_0$ .

We will again work out the likelihood ratio test. We first write down the “conditional” likelihood of the data:

$$K((X_n)_{n=1}^N | (c_n)_{n=1}^N, \theta) = \prod_{n=1}^N \theta_{X_n|c_n} = \prod_{x \in \mathcal{X}} \prod_{c \in \mathcal{C}} \theta_{x|c}^{N_{x|c}}$$

where we used the counts as a sufficient statistic of the data. Maximizing the likelihood with respect to the parameters, we obtain the maximum likelihood estimator

$$\hat{\theta}_{x|c} = \frac{N_{x|c}}{N_{+|c}},$$

i.e., the fractions of the different outcomes within each subgroup with  $C = c$ . Under  $H_0$ , we can write  $\theta_{x|c} = \theta_{x|*}$  for some  $\theta_{x|*} \in \Theta_0$ , and the likelihood simplifies:

$$\theta \in \Theta_0 \implies \prod_{x \in \mathcal{X}} \prod_{c \in \mathcal{C}} \theta_{x|c}^{N_{x|c}} = \prod_{x \in \mathcal{X}} \theta_{x|*}^{N_{x|+}}.$$



The maximum likelihood estimator under  $H_0$  is then

$$\hat{\theta}_{x|c}^0 = \frac{N_{x|+}}{N}.$$

The likelihood ratio is obtained by dividing the likelihood for  $\hat{\theta}$  by the likelihood for  $\hat{\theta}^0$ :

$$\frac{\sup_{\theta \in \Theta_0 \cup \Theta_1} K((X_n)_{n=1}^N | (c_n)_{n=1}^N, \theta)}{\sup_{\theta \in \Theta_0} K((X_n)_{n=1}^N | (c_n)_{n=1}^N, \theta)} = \prod_{x \in \mathcal{X}} \prod_{c \in \mathcal{C}} \left( \frac{\hat{\theta}_{x|c}}{\hat{\theta}_{x|c}^0} \right)^{N_{x|c}} = \prod_{x \in \mathcal{X}} \prod_{c \in \mathcal{C}} \left( \frac{N_{x|c}}{N_{+|c}} \frac{N}{N_{x|+}} \right)^{N_{x|c}}. \quad (64)$$

This is of *exactly* the same form as the likelihood ratio (58) for testing  $X \perp\!\!\!\perp_{P(X,C)} C$  with two random variables  $X, C$ . The likelihood ratio test statistic is defined as 2 times the natural logarithm of this ratio:

$$G_N := 2 \log \frac{\sup_{\theta \in \Theta_0 \cup \Theta_1} K((X_n)_{n=1}^N | (c_n)_{n=1}^N, \theta)}{\sup_{\theta \in \Theta_0} K((X_n)_{n=1}^N | (c_n)_{n=1}^N, \theta)} = 2 \sum_{x \in \mathcal{X}} \sum_{c \in \mathcal{C}} N_{x|c} \log \frac{N_{x|c} N}{N_{x|+} N_{+|c}}.$$

We (miraculously?!) arrived at the same test statistic as before. This time, we cannot make use of the general result on the asymptotic distribution under the null hypothesis of likelihood ratio tests. Indeed, that result pertains when dealing with a likelihood with only a finite number of parameters to be estimated, whereas here we have an asymptotically infinite number of “parameters”  $c_1, c_2, \dots$ . Therefore, we will resort to a more direct analysis of the case at hand. The end result will recover the previous results for random variables as a special case. Perhaps surprisingly, it turns out that under reasonable assumptions on the sequence  $c_1, c_2, \dots$  we can apply the standard  $G$  test and ignore the non-random nature of the  $c_n$ 's.

We will start with rewriting the test statistic. We introduce the space

$$\Theta_C := \left\{ \gamma \in \prod_{c \in \mathcal{C}} [0, 1] : \sum_{c \in \mathcal{C}} \gamma_c = 1 \right\}.$$

Consider now the function

$$\begin{aligned} g : \Theta \times \Theta_C : (\theta, \gamma) &\mapsto := \sum_{c \in \mathcal{C}} \gamma_c \sum_{x \in \mathcal{X}} \theta_{x|c} \log \frac{\theta_{x|c}}{\sum_{c \in \mathcal{C}} \gamma_c \theta_{x|c}} \\ &= \sum_{c \in \mathcal{C}} \gamma_c \text{KL} \left( \theta_{X|c} \parallel \sum_{c \in \mathcal{C}} \gamma_c \theta_{X|c} \right), \end{aligned} \quad (65)$$

where we introduced the notation  $\theta_{X|c} := (\theta_{x|c})_{x \in \mathcal{X}} \in \mathbb{R}^{\mathcal{X}}$ , and where the Kullback-Leibler divergence between two probability distributions  $P, Q \in \mathcal{P}(\mathcal{X})$  is defined as:

$$\text{KL}(P \parallel Q) := \sum_{x \in \mathcal{X}} P(X = x) \log \frac{P(X = x)}{Q(X = x)}$$

(and we identified a probability distribution on  $\mathcal{X}$  with its probability mass function, encoded as a parameter vector in  $\mathbb{R}^{\mathcal{X}}$ ). Note that

$$G_N = 2Ng(\hat{\theta}, \hat{\gamma})$$

with

$$\hat{\gamma}_c := \frac{N_{+|c}}{N}$$

for all  $c \in \mathcal{C}$ .<sup>54</sup> The components of  $\hat{\gamma}$  are just the fractions of observations with a certain value of  $c$ . The asymptotic analysis will be conditional on the sequence of  $\hat{\gamma}$ 's, or equivalently, on the sequence of counts  $N_{+|c}$ .

Because the counts  $(N_{x|c})_{x \in \mathcal{X}}$  have a multinomial distribution for each  $c \in \mathcal{C}$ , we can calculate that

$$\mathbb{E}\hat{\theta}_{x|c} = \frac{1}{N_{+|c}} \sum_{n=1}^N \mathbb{1}_c(c_n) \theta_{x|c} = \theta_{x|c}$$

and

$$\begin{aligned} \text{Cov}(\hat{\theta}_{x|c}, \hat{\theta}_{x'|c'}) &= \frac{1}{N_{+|c}N_{+|c'}} \sum_{n=1}^N \sum_{n'=1}^N \text{Cov}(\mathbb{1}_{(x,c)}(X_n, c_n) \mathbb{1}_{(x',c')}(X_{n'}, c_{n'})) \\ &= \frac{1}{N_{+|c}N_{+|c'}} \sum_{n=1}^N \mathbb{1}_c(c_n) \mathbb{1}_{c'}(c_n) (\delta_{xx'} \theta_{x|c} - \theta_{x|c} \theta_{x'|c'}) \\ &= \frac{1}{N_{+|c}} \delta_{cc'} \theta_{x|c} (\delta_{xx'} - \theta_{x'|c'}). \end{aligned}$$

Assume that  $N_{+|c} \rightarrow \infty$  for every  $c \in \mathcal{C}$ . Because

$$N_{x|c} = \sum_{\substack{n=1 \\ c_n=c}}^N \mathbb{1}_x(X_n),$$

where the  $\mathbb{1}_x(X_n)$  can be seen as i.i.d. vectors in  $\mathbb{R}^{\mathcal{X}}$ , we can apply the multivariate central limit theorem to conclude that

$$\sqrt{N_{+|c}}(\hat{\theta}_{x|c} - \theta_{x|c}) \rightsquigarrow \mathcal{N}(0, \Sigma) \tag{66}$$

with covariance matrix  $\Sigma$  with entries

$$(\Sigma)_{x|c, x'|c'} = \delta_{cc'} \theta_{x|c} (\delta_{xx'} - \theta_{x'|c'}). \tag{67}$$

The central limit theorem thus provides the rate at which the ML estimate  $\hat{\theta}$  converges to the true parameter  $\theta$ . The likelihood ratio statistic  $G_N$  is a function of  $\hat{\theta}$  and  $\hat{\gamma}$ . The asymptotic analysis of  $G_N$  can be obtained by performing a Taylor expansion of  $g$  around the true  $\theta$ . This expansion can be done “uniformly” in  $\hat{\gamma}$ .

---

<sup>54</sup>We write  $\hat{\gamma}$  rather than  $\gamma$  to indicate that it is an ( $N$ -dependent) function of the observed data, rather than a “true” fixed quantity.

**Lemma 10.3.1.** Let  $\theta \in \Theta_0$  be positive, i.e., such that there exists a  $\delta > 0$  with  $\theta_{x|c} \geq \delta$  for all  $x \in \mathcal{X}, c \in \mathcal{C}$ . Let  $\gamma \in \Theta_C$ . For  $\hat{\theta}$  with  $\|\hat{\theta} - \theta\|$  sufficiently small and  $\hat{\theta} \xrightarrow{P} \theta$ ,

$$g(\hat{\theta}, \gamma) = \frac{1}{2}(\hat{\theta} - \theta)^T \nabla_{\theta}^2 g(\theta, \gamma)(\hat{\theta} - \theta) + o_P(\|\hat{\theta} - \theta\|^2), \quad (68)$$

with

$$\nabla_{\theta}^2 g(\theta, \gamma) = \text{diag}\left(\frac{1}{\theta_{X|*}}\right) \otimes \left(\text{diag}(\sqrt{\gamma})(I_C - \sqrt{\gamma}\sqrt{\gamma}^T)\text{diag}(\sqrt{\gamma})\right)$$

where  $\theta_{X|*} = (\theta_{x|c})_{x \in \mathcal{X}}$  for an arbitrary  $c \in \mathcal{C}$ .<sup>55</sup> The remainder term  $o_P(\|\hat{\theta} - \theta\|^2)$  can be chosen to be a function  $f(\theta, \hat{\theta})$  that does not depend on  $\gamma$ , and for which

$$\frac{f(\theta, \hat{\theta})}{\|\theta - \hat{\theta}\|^2} \xrightarrow{P} 0.$$

*Proof.* The second order Taylor expansion of  $g$  around the true  $\theta$  is:

$$g(\theta + \epsilon, \gamma) = g(\theta, \gamma) + \epsilon^T \nabla_{\theta} g(\theta, \gamma) + \frac{1}{2} \epsilon^T \nabla_{\theta}^2 g(\theta, \gamma) \epsilon + o(\|\epsilon\|^2),$$

where  $\nabla_{\theta} g$  is the (partial) gradient of  $g$  with respect to  $\theta$ , and  $\nabla_{\theta}^2 g$  is the Hessian of  $g$  with respect to  $\theta$ , and the remainder term  $o(\|\epsilon\|^2)$  can be taken of the form  $M\|\epsilon\|^3$  if the third-order partial derivatives of  $g$  at  $(\theta, \gamma)$  are bounded. For a random  $\epsilon = \hat{\theta} - \theta$  we obtain

$$g(\hat{\theta}, \gamma) - g(\theta, \gamma) = (\hat{\theta} - \theta)^T \nabla_{\theta} g(\theta, \gamma) + \frac{1}{2}(\hat{\theta} - \theta)^T \nabla_{\theta}^2 g(\theta, \gamma)(\hat{\theta} - \theta) + o_P(\|\hat{\theta} - \theta\|^2), \quad (69)$$

where the remainder term is now random, and converges in probability to 0 at rate  $\|\hat{\theta} - \theta\|^2$ . We will proceed by calculating the terms in the Taylor expansion.

The Kullback-Leibler divergence has the important property that  $\text{KL}(P \parallel Q) \geq 0$  and  $\text{KL}(P \parallel Q) = 0 \iff P = Q$  for all  $P, Q \in \mathcal{P}(\mathcal{X})$ . Together with the definition (65), this immediately implies that under  $H_0$ ,  $g(\theta, \gamma) = 0$  for all  $\gamma \in \Theta_C$ .

The gradient  $\nabla_{\theta} g$  of  $g$  w.r.t.  $\theta$  has components:

$$\frac{\partial g}{\partial \theta_{x|c}} = \gamma_c \log \frac{\theta_{x|c}}{\sum_{c' \in \mathcal{C}} \gamma_{c'} \theta_{x|c'}}.$$

Under  $H_0$ ,  $\theta_{x|c} = \theta_{x|*}$  for all  $x \in \mathcal{X}, c \in \mathcal{C}$ , and it follows that the gradient vanishes for all  $\gamma \in \Theta_C$ .

Next, the Hessian w.r.t.  $\theta$ :

$$\frac{\partial^2 g}{\partial \theta_{x|c} \partial \theta_{x'|c'}} = \gamma_c \left( \frac{1}{\theta_{x|c}} \delta_{xx'} \delta_{cc'} - \frac{\gamma_{c'}}{\sum_{c'' \in \mathcal{C}} \gamma_{c''} \theta_{x|c''}} \delta_{xx'} \right).$$

<sup>55</sup>Here, we used the Kronecker product notation for matrices, and  $\text{diag}(v)$  is a diagonal matrix with the components of vector  $v$  on the diagonal.

Under  $H_0$ ,  $\theta_{x|c} = \theta_{x|*}$  for all  $x \in \mathcal{X}, c \in \mathcal{C}$ , this simplifies to

$$\frac{\partial^2 g}{\partial \theta_{x|c} \partial \theta_{x'|c'}} = \frac{1}{\theta_{x|*}} \gamma_c (\delta_{xx'} \delta_{cc'} - \gamma_{c'} \delta_{xx'}).$$

By using the Kronecker product notation, this can be written as stated in the lemma.

Finally, to obtain the remainder term, we calculate the third order partial derivatives:

$$\frac{\partial^3 g}{\partial \theta_{x|c} \partial \theta_{x'|c'} \partial \theta_{x''|c''}} = \gamma_c \delta_{x,x''} \delta_{x,x'} \left( -\frac{1}{\theta_{x|c}^2} \delta_{c,c''} \delta_{cc'} + \frac{\gamma_{c'} \gamma_{c''}}{(\sum_{c''' \in \mathcal{C}} \gamma_{c'''} \theta_{x|c'''})^2} \right).$$

These can be bounded uniformly in  $\gamma$ , using the assumption that all components of  $\theta$  are bounded away from zero.  $\square$

We are now ready to prove the following result on the asymptotic distribution of  $G_N$  under the null hypothesis, which is (surprisingly?) similar to Proposition 10.1.1.

**Proposition 10.3.2.** *Let  $\theta \in \Theta$  be positive, i.e., such that  $\theta_{x|c} > 0$  for all  $x \in \mathcal{X}, c \in \mathcal{C}$ . Assume that  $N_{+|c} \rightarrow \infty$  for all  $c \in \mathcal{C}$ . Under  $H_0 : X \perp_{K(X|C)} C$ , the likelihood ratio test statistic (64) converges to a  $\chi^2$  distribution,*

$$G_N \rightsquigarrow \chi_\nu^2,$$

with  $\nu = (|\mathcal{X}| - 1)(|\mathcal{C}| - 1)$  degrees of freedom.

*Proof.* We first note that  $\hat{\theta} \xrightarrow{a.s.} \theta$  (for any  $\theta \in \Theta$ ) by applying the strong law of large numbers. Indeed, since the  $X_n$  are assumed to be conditionally i.i.d. given  $c_n$ , and

$$\mathbb{E}(\mathbb{1}_{(x,c)}(X_n, c_n)) = \theta_{x|c} \mathbb{1}_c(c_n)$$

for all  $x \in \mathcal{X}, c \in \mathcal{C}$ , and  $N_{+|c} \rightarrow \infty$  for all  $c \in \mathcal{C}$ , we conclude that  $N_{x|c}/N_{+|c} \xrightarrow{a.s.} \theta_{x|c}$  for all  $x \in \mathcal{X}, c \in \mathcal{C}$  by the strong law of large numbers.

Since this implies convergence in probability  $\hat{\theta} \xrightarrow{P} \theta$ , Lemma 10.3.1 gives that under  $H_0$ ,

$$g(\hat{\theta}, \hat{\gamma}) = \frac{1}{2} (\hat{\theta} - \theta)^T \nabla_{\hat{\theta}}^2 g(\theta, \hat{\gamma}) (\hat{\theta} - \theta) + o_P(\|\hat{\theta} - \theta\|^2)$$

where

$$\nabla_{\hat{\theta}}^2 g(\theta, \hat{\gamma}) = \text{diag} \left( \frac{1}{\theta_{X|*}} \right) \otimes \left( \text{diag}(\sqrt{\hat{\gamma}}) (I_C - \sqrt{\hat{\gamma}} \sqrt{\hat{\gamma}}^T) \text{diag}(\sqrt{\hat{\gamma}}) \right)$$

and the remainder term does not depend on  $\hat{\gamma}$ . Defining

$$S_N := \sqrt{N} \text{diag} \left( \frac{1}{\sqrt{\theta_{X|*}}} \otimes \sqrt{\hat{\gamma}} \right) (\hat{\theta} - \theta)$$

where  $\otimes$  denotes the Kronecker product of two vectors (in this case the all-ones vector  $\mathbf{1}_{\mathcal{X}} \in \mathbb{R}^{\mathcal{X}}$  and the vector  $\sqrt{\hat{\gamma}} \in \mathbb{R}^{\mathcal{C}}$ ), and the orthogonal projection<sup>56</sup>

$$\hat{\Gamma} := I_{\mathcal{X}} \otimes \left( I_C - \sqrt{\hat{\gamma}} \sqrt{\hat{\gamma}}^T \right),$$

<sup>56</sup>That is,  $\hat{\Gamma}^T = \hat{\Gamma} = \hat{\Gamma}^2$ .

we can write

$$N(\hat{\theta} - \theta)^T \nabla_{\theta}^2(\theta, \hat{\gamma})(\hat{\theta} - \theta) = S_N^T \hat{\Gamma} S_N = \|\hat{\Gamma} S_N\|^2.$$

Under  $H_0$ ,  $\theta_{x|c} = \theta_{x|*}$  for all  $c \in \mathcal{C}$ , and therefore (66) simplifies to:

$$\sqrt{N} \sqrt{\hat{\gamma}_c} (\hat{\theta}_{x|c} - \hat{\theta}_{x|*}) \rightsquigarrow \mathcal{N}(0, \Sigma)$$

where the covariance matrix  $\Sigma$  from (67) simplifies to

$$(\Sigma)_{x|c, x'|c'} = \delta_{cc'} \theta_{x|*} (\delta_{xx'} - \theta_{x'|*}).$$

Hence, scaling with  $\sqrt{\theta_{x|*}}$  gives

$$S_N = \sqrt{N} \sqrt{\hat{\gamma}_c} \frac{1}{\sqrt{\theta_{x|*}}} (\hat{\theta}_{x|c} - \theta_{x|*}) \rightsquigarrow \mathcal{N}(0, \tilde{\Sigma})$$

where the covariance matrix  $\tilde{\Sigma}$  from (67) simplifies to

$$(\tilde{\Sigma})_{x|c, x'|c'} = \delta_{cc'} \left( \delta_{xx'} - \sqrt{\theta_{x|*}} \sqrt{\theta_{x'|*}} \right),$$

which can be written using Kronecker product notation as

$$\tilde{\Sigma} = \left( I_{\mathcal{X}} - \sqrt{\theta_{X|*}} \sqrt{\theta_{X|*}^T} \right) \otimes I_{\mathcal{C}}.$$

Let  $\hat{V} \in SO(\mathbb{R}^c)$  be rotations that map  $\sqrt{\hat{\gamma}}$  to  $e_1$ . Then

$$I_{\mathcal{C}} - \sqrt{\hat{\gamma}} \sqrt{\hat{\gamma}}^T = P_{\sqrt{\hat{\gamma}}^\perp} = \hat{V}^T P_{e_1^\perp} \hat{V}$$

(where  $P_{v^\perp}$  is the orthogonal projection on the subspace orthogonal to  $v$ ), and therefore

$$\|\hat{\Gamma} S_N\|^2 = \|(I_{\mathcal{X}} \otimes \hat{V}^T P_{e_1^\perp} \hat{V}) S_N\|^2 = \|(I_{\mathcal{X}} \otimes P_{e_1^\perp})(I_{\mathcal{X}} \otimes \hat{V}) S_N\|^2.$$

We can apply Lemma 10.3.3 to conclude that  $(I_{\mathcal{X}} \otimes \hat{V}) S_N \rightsquigarrow \mathcal{N}(0, \tilde{\Sigma})$  as well (even though  $\hat{V}$  is an  $N$ -dependent rotation). Then

$$(I_{\mathcal{X}} \otimes P_{e_1^\perp})(I_{\mathcal{X}} \otimes \hat{V}) S_N \rightsquigarrow \mathcal{N}\left(0, P_{\sqrt{\theta_{X|*}}^\perp} \otimes P_{e_1^\perp}\right)$$

and hence

$$\|\hat{\Gamma} S_N\|^2 \rightsquigarrow \chi_\nu^2$$

with  $\nu = (|\mathcal{X}| - 1)(|\mathcal{C}| - 1)$ . With Slutsky's Lemma, also

$$G_N = 2Ng(\hat{\theta}, \hat{\gamma}) = N(\hat{\theta} - \theta)^T \nabla_{\theta}^2(\theta, \hat{\gamma})(\hat{\theta} - \theta) + N o_P(\|\hat{\theta} - \theta\|^2) \rightsquigarrow \chi_\nu^2.$$

□

So, also the asymptotic distribution under the null hypothesis is the same as for the case of two random variables, even though we used a different parameterization, and we relaxed the assumption that the  $c_n$ 's are i.i.d.: we only assumed that  $N_{+|c} \rightarrow \infty$  for each  $c$ .<sup>57</sup> In particular, this result applies also to the case when  $C$  is a random variable. Hence, we have reobtained Proposition 10.1.1 as a special case.

**Lemma 10.3.3.** *Let  $Q$  be a rotationally symmetric probability measure on the standard Borel space  $\mathbb{R}^k$  (i.e.,  $Q \circ U = Q$  for all  $U \in SO(\mathbb{R}^k)$ ), and  $P_1, P_2, \dots$  a sequence of probability measures on  $\mathbb{R}^k$ . Then*

$$P_n \rightsquigarrow Q \iff P_n \circ U_n \rightsquigarrow Q$$

for any sequence  $U_1, U_2, \dots$  of rotations in  $SO(\mathbb{R}^k)$ .

*Proof.* We make use of the Lévy-Prokhorov metrization of the weak topology. For two probability distributions  $P, Q$  on  $\mathbb{R}^k$  (with its Borel  $\sigma$ -algebra  $\mathcal{B}_{\mathbb{R}^k}$ ), it is defined as

$$d(P, Q) := \inf\{\epsilon > 0 : \forall A \in \mathcal{B}_{\mathbb{R}^k} : P(A) \leq Q(A^\epsilon) + \epsilon \wedge Q(A) \leq P(A^\epsilon) + \epsilon\},$$

where  $A^\epsilon := \cup_{a \in A} B_\epsilon(a)$  with  $B_\epsilon(a) := \{v \in \mathbb{R}^k : \|v - a\| < \epsilon\}$ . The Lévy-Prokhorov metric is invariant under rotations. Indeed, for  $U \in SO(\mathbb{R}^k)$ , we have that  $UA \in \mathcal{B}_{\mathbb{R}^k} \iff A \in \mathcal{B}_{\mathbb{R}^k}$ , and  $(UA)^\epsilon = U(A^\epsilon)$  for  $A \in \mathcal{B}_{\mathbb{R}^k}$  and  $\epsilon > 0$ , hence  $d(P \circ U, Q \circ U) = d(P, Q)$  for all probability measures  $P, Q$  on  $\mathbb{R}^k$ . The rotation invariance of the Lévy-Prokhorov metric implies that if  $Q$  is rotationally symmetric (i.e.,  $Q = Q \circ U$  for all  $U \in SO(\mathbb{R}^n)$ ), then

$$\begin{aligned} P_n \rightsquigarrow Q &\iff d(P_n, Q) \rightarrow 0 \\ &\iff d(P_n \circ U_n, Q \circ U_n) \rightarrow 0 \\ &\iff d(P_n \circ U_n, Q) \rightarrow 0 \\ &\iff P_n \circ U_n \rightsquigarrow Q \end{aligned}$$

□

Summarizing, we started from quite a different sampling scheme, and did not treat  $C$  as a random variable, yet we ended up with exactly the same likelihood ratio test for testing  $X \perp_{K(X|C)} C$  as we derived for testing the independence  $X \perp_{P(X,C)} C$  between two random variables.

What about the consistency of this test? To show consistency, it turns out that we need a stronger assumption than  $N_{+|c} \rightarrow \infty$ , but it will still be weaker than the i.i.d. assumption we made for the case that  $C$  is random.

**Corollary 10.3.4.** *Consider an infinite sequence of  $G$  tests performed on the first  $N$  samples of an infinitely large data set  $(X_n, c_n)_{n=1}^\infty$ , where one accepts  $H_1 : X \not\perp C$  if  $G_N \geq \tau_N$ , and otherwise accepts  $H_0 : X \perp C$ , for some given sequence of thresholds  $\tau_N$ .*

<sup>57</sup>If some of the  $N_{+|c}$  stay finite asymptotically, then these  $c$ 's can be ignored, and we still get asymptotically a chi-square distribution, but with less degrees of freedom.

Assume that  $\theta \in \Theta$  is positive, i.e., such that  $\theta_{x|c} > 0$  for all  $x \in \mathcal{X}, c \in \mathcal{C}$ . Assume further that the fractions  $N_{+|c}/N \rightarrow \infty$  are bounded away from zero asymptotically, i.e., there exists  $\epsilon > 0$  such that for all  $c \in \mathcal{C}$ ,  $N_{+|c}/N \geq \epsilon$  for large  $N$ . Then this sequence of tests is asymptotically consistent if  $\tau_N \rightarrow \infty$  but  $\tau_N/N \rightarrow 0$ .

*Proof.* In the proof of Proposition 10.3.2 we already saw that  $\hat{\theta} \xrightarrow{\text{a.s.}} \theta$  for any  $\theta \in \Theta$ .

Assume that  $H_1$  holds, i.e.,  $\theta \in \Theta_1$ . Then  $g(\theta, \gamma) = I(\gamma, \theta \circ \gamma) > 0$  for all  $\gamma \in \Theta_C$ .  $\Theta_1$  is open in  $\Theta$ , so  $\bar{B}_\delta(\theta) \cap \Theta = \{\tilde{\theta} \in \Theta : \|\theta - \tilde{\theta}\| \leq \delta\} \subseteq \Theta_1$  for  $\delta$  small enough. Since  $\hat{\theta} \xrightarrow{\text{a.s.}} \theta$ ,  $\hat{\theta} \in \bar{B}_\delta(\theta) \cap \Theta$  for large  $N$  a.s.. By assumption, for large  $N$   $\hat{\gamma} \in \{\tilde{\gamma} \in [\epsilon, 1]^{|\mathcal{C}|} : \tilde{\gamma}_+ = 1\}$ , which is a closed subset of  $\Theta_C$ . Since  $I$  is continuous, it attains a (positive) minimum value over the closed subset  $(\bar{B}_\delta(\theta) \cap \Theta) \times \{\tilde{\gamma} \in [\epsilon, 1]^{|\mathcal{C}|} : \tilde{\gamma}_+ = 1\} \subseteq \Theta \times \Theta_C$ . Hence,  $G_N/N = 2g(\hat{\theta}, \hat{\gamma})$  a.s. has a positive lower bound for large  $N$ . We can now reason analogously as in the proof of Corollary 10.1.3 to conclude that the probability of a Type II error vanishes asymptotically.

For an asymptotic estimate of the probability of a Type I error, we can make use of Proposition 10.3.2, which states that  $G_N \rightsquigarrow \chi_\nu^2$  under  $H_0$ . This part of the proof is identical to the corresponding part of the proof of Corollary 10.1.3.  $\square$

The conditions in this corollary are sufficient, but not necessary. For example, not all the rates  $N_{+|c}$  have to be lower bounded, it suffices if this is the case for a subset of  $\mathcal{C}$  for which the distributions  $K(X | C = c)$  differ. It also shows how consistency could fail: e.g., if the distributions  $K(X | C = c)$  only differ on some subset of  $\mathcal{C}$ , but that subset is not observed sufficiently often asymptotically. This is in line with the intuition that when testing for the presence of a causal effect of  $C$  on  $X$  in a controlled setting (not necessarily randomized, i.e., as in Proposition 9.5.1), if nothing is known about how  $X$  might depend on  $C$ , it is best to gather sufficient data for each value that  $C$  can take.

## 10.4. The general categorical case

Finally, let us consider the most general case of testing a conditional independence involving transitional random variables (including “purely random” and “purely non-random” variables as special cases). Again, we will restrict ourselves to the case that all variables take values in finite spaces. We will formulate a general version of the  $G$  test.

Suppose we have three transitional random variables  $X, Y, Z$  and an input variable  $C$ , taking values in spaces  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}, \mathcal{C}$ , respectively. Assume that all the spaces  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}, \mathcal{C}$  are finite. Suppose we have a kernel  $K(X, Y, Z | C) : \mathcal{C} \dashrightarrow \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$ . We formulate a statistical test for testing the conditional independence

$$H_0 : X \underset{K(X,Y,Z|C)}{\perp\!\!\!\perp} Y | Z$$

against the alternative

$$H_1 : X \not\underset{K(X,Y,Z|C)}{\perp\!\!\!\perp} Y | Z.$$

The null hypothesis is equivalent, by definition, to the existence of a kernel  $K(X|Z)$  such that

$$K(X, Y, Z | C) = K(X | Z) \otimes K(Y, Z | C).$$

It will turn out to be helpful to consider the equivalent hypotheses

$$H_0 : X \perp\!\!\!\perp_{K(X,Y,Z|C)} Y, C | Z$$

against the alternative

$$H_1 : X \not\perp\!\!\!\perp_{K(X,Y,Z|C)} Y, C | Z$$

instead.

Suppose we obtain data  $(X_n, Y_n, Z_n, c_n)_{n=1}^N$  such that the  $(X_n, Y_n, Z_n)$  are conditionally independent and identically distributed given  $c_n$  for all  $n = 1, \dots, N$ . In other words, we assume the data is sampled from the following Markov kernel:

$$K((X_n, Y_n, Z_n)_{n=1}^N | (c_n)_{n=1}^N) = \bigotimes_{n=1}^N K(X_n, Y_n, Z_n | C_n = c_n),$$

where each  $K(X_n, Y_n, Z_n | C_n)$  is a copy of the (“true” but unknown kernel)  $K(X, Y, Z | C)$ . We parameterize this kernel  $K(X, Y, Z | C)$  as

$$K(X = x, Y = y, Z = z | C = c) = \theta_{xyz|c}$$

with  $\theta$  in

$$\Theta = \left\{ \theta \in \prod_{(x,y,z,c) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z} \times \mathcal{C}} : \theta_{+++|c} = 1 \ \forall c \in \mathcal{C} \right\}.$$

We will not work out the details, as these are analogous to what we have seen before, but will directly formulate the likelihood ratio test statistic:

$$G_N = 2 \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}} \sum_{c \in \mathcal{C}} N_{xyzc} \log \frac{N_{xyzc} N_{+++|c}}{N_{x+z|c} N_{+yz|c}}. \quad (70)$$

With the correspondence  $(X, (Y, C), Z) \leftrightarrow (X, Y, Z)$ , this likelihood ratio statistic is seen to be identical to (61), the one for the case of three random variables. Note that this likelihood ratio treats  $C$  and  $Y$  on equal footing. This, again, suggests that for the asymptotic analysis we will obtain a similar result as that for the conditional  $G$  test with purely random variables, albeit under milder assumptions on the sampling scheme for  $C$ . We will not work out the details here, but directly formulate the result.

where we used the counts as a sufficient statistic of the data.

**Proposition 10.4.1.** *Let  $\theta \in \Theta$  be positive, i.e., such that  $\theta_{xyz|c} > 0$  for all  $x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}, c \in \mathcal{C}$ . Assume that  $N_{+++|c} \rightarrow \infty$  for all  $c \in \mathcal{C}$ . Under  $H_0 : X \perp\!\!\!\perp_{K(X,Y,Z|C)} Y, C | Z$ , the likelihood ratio test statistic (70) converges to a  $\chi^2$  distribution,*

$$G_N \rightsquigarrow \chi_\nu^2,$$

with  $\nu = |\mathcal{Z}|(|\mathcal{X}| - 1)(|\mathcal{Y}||\mathcal{C}| - 1)$  degrees of freedom.



*Proof.* Note that the likelihood ratio test statistic (70) is a sum over  $z \in \mathcal{Z}$  of a likelihood ratio test statistic of the form (64) (where  $(Y, C)$  in the former corresponds with  $C$  in the latter).  $\square$

We also obtain a similar result as before on the asymptotic consistency.

**Corollary 10.4.2.** *Consider an infinite sequence of  $G$  tests performed on the first  $N$  samples of an infinitely large data set  $(X_n, Y_n, Z_n, c_n)_{n=1}^{\infty}$ , where one decides*

$$\begin{cases} H_0 : X \perp\!\!\!\perp_{K(X,Y,Z|C)} Y, C | Z & \text{if } G_N < \tau_N, \\ H_1 : X \not\perp\!\!\!\perp_{K(X,Y,Z|C)} Y, C | Z & \text{if } G_N \geq \tau_N, \end{cases}$$

for some given sequence of thresholds  $\tau_N$ . Assume that  $\theta \in \Theta$  is positive, i.e., such that  $\theta_{xyz|c} > 0$  for all  $x \in \mathcal{X}, y \in \mathcal{Y}, z \in \mathcal{Z}, c \in \mathcal{C}$ . Assume further that the fractions  $N_{+++|c}/N \rightarrow \infty$  are bounded away from zero asymptotically, i.e., there exists  $\epsilon > 0$  such that for all  $c \in \mathcal{C}$ ,  $N_{+++|c}/N > \epsilon$  for large  $N$ . Then this sequence of tests is asymptotically consistent if  $\tau_N \rightarrow \infty$  but  $\tau_N/N \rightarrow 0$ .

## 11. The Fast Causal Inference Algorithm

In this final chapter, we present an extension of the Fast Causal Inference (FCI) algorithm. The FCI algorithm is one of the highlights in the field of constraint-based causal discovery. It was originally designed for purely observational (non-experimental) data and relied on the assumption of acyclicity, but allowed for latent variables (either marginalized over, or conditioned upon in case of selection bias) [SMR95, SMR99]. Later, the algorithm was augmented with additional ‘orientation rules’ and this augmented FCI algorithm was shown to be complete in a certain sense [Zha08].<sup>58</sup> While the FCI algorithm was originally designed for the acyclic setting, it was recently discovered that it also works in case cycles are present (more specifically, for  $\sigma$ -faithful simple SCMs) [MC20]. While that work made a simplification by assuming no selection bias, we prove here that the FCI algorithm is sound even when the data is generated according to a conditional Markov kernel induced by a  $\sigma$ -faithful simple SCM, or in other words, that it can cope with the possible presence of cycles *and* selection bias. In other recent work, the FCI algorithm has been extended to incorporate prior knowledge regarding exogeneity and unconfoundedness of context variables (of the same type that we used for modeling the treatment variable in randomized controlled trials) [MMC20]. Such an extension with exogenous input nodes allows to generalize the idea of causal discovery in a randomized controlled trial to multiple treatment and outcome variables. We also incorporate that extension here, and provide a ‘unified’ extended FCI algorithm that allows for cycles (under certain assumptions), selection bias and exogenous input variables.

### 11.1. Modeling selection bias

We model selection bias as follows.

**Notation 11.1.1.** *We will assume the existence of a simple SCM  $M = (J, V^+, W, \mathcal{X}, P, f)$ , with endogenous variables  $V^+ = V \dot{\cup} S$ .<sup>59</sup> We assume that only the endogenous variables in  $V$  are observed, as well as the exogenous input variables  $J$ . The latent variables in  $S$  have the role of latent selection variables. In other words, we will assume that the data is distributed according to the marginal Markov kernel of  $M$  on  $V$  after conditioning on  $S$ :*

$$P_M(X_V \mid X_S \in \xi_S, \text{do}(X_J))$$

*for some (unknown) measurable subset  $\xi_S \subseteq \mathcal{X}_S$ .<sup>60</sup>*

---

<sup>58</sup>The augmented version of [Zha08] is often referred to simply as ‘the FCI algorithm’, and we will do so here as well.

<sup>59</sup>Alternatively, one could assume endogenous variables  $V^+ = V \dot{\cup} S \dot{\cup} L$  where  $S$  are used as selection variables. We could then start by marginalizing out the variables in  $L$  to arrive at the setting that is our starting point here.

<sup>60</sup>Note that we did not include the exogenous random variables here; we are treating them as if they were unobserved. If some (or all) exogenous random variables are observed, then we have three options: (i) include them, but ignore the additional background knowledge that they are independent (that

This models a process that *filters* the data according to the value of  $X_S$  (like in a rejection sampler).<sup>61</sup> One practical example of such a filtering process leading to selection bias is the following.

**Example 11.1.2.** *After the exam, the teacher hands out evaluation forms to the students and asks them to evaluate the course. When analyzing this evaluation, the teacher needs to be aware of possible selection bias. For example, if the students that were most dissatisfied with the course already dropped out earlier on and never filled in the evaluation form, the observed evaluations may not be representative of the opinions of all the students that started the course. Here, the selection variable may be the binary indicator ‘student filled in the evaluation form’.*

Our goal in the rest of this chapter will be to deduce as much as possible about the causal graph  $G_{V+|J}(M)$  from the marginal Markov kernel  $P_M(X_V | X_S \in \xi_S, \text{do}(X_J))$ .

## 11.2. Inducing walks

The key notion in this chapter is that of  $\sigma$ -*inducing walks*. We define  $\sigma$ -inducing walk as a generalization to CDMGs of the notion of inducing path [VP90] in DAGs.

**Definition 11.2.1.** *Let  $G$  be a CDMG with input nodes  $J$ , output nodes  $V^+ = V \dot{\cup} S$  and let  $i, j \in V \dot{\cup} J$  be distinct nodes. A walk  $\pi$  in  $G$  between  $i$  and  $j$  is called  $\sigma$ -inducing given  $S$  if each collider on  $\pi$  is in  $\text{Anc}_G(\{i, j\} \cup S)$ , and each non-endpoint non-collider on  $\pi$  is unblockable. If it is a path, it is called a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$ .<sup>62</sup>*

If two nodes are adjacent in  $G$ , any edge connecting the two is a  $\sigma$ -inducing walk (path) given  $S$  between them, for any  $S$ . Figure 26 shows some simple nontrivial examples of  $\sigma$ -inducing paths.

**Lemma 11.2.2.** *Let  $G$  be a CDMG with input nodes  $J$ , output nodes  $V^+ = V \dot{\cup} S$  and let  $i, j \in V \dot{\cup} J$  be distinct nodes. If  $i \in \text{Sc}_G(j)$  then there exists a  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  and  $j$ .*

*Proof.* There exists a directed path in  $G$  from  $i$  to  $j$  that is entirely contained in  $\text{Sc}_G(j)$ , and therefore all its non-endpoint nodes are unblockable non-colliders.  $\square$

The notion of  $\sigma$ -inducing walk has the following important properties.

---

is, treating them as if they were endogenous), (ii) make endogenous copies and include those; (iii) include them and make use of the additional background knowledge that they are independent. Here, we chose for the second option.

<sup>61</sup>It is helpful to think about such a filtering step as an intervention on the population level, but may lead to confusion when interpreted as an intervention on the individual level.

<sup>62</sup>In most of the literature, the graph is assumed to be acyclic, and then the notion is referred to simply as “inducing”.

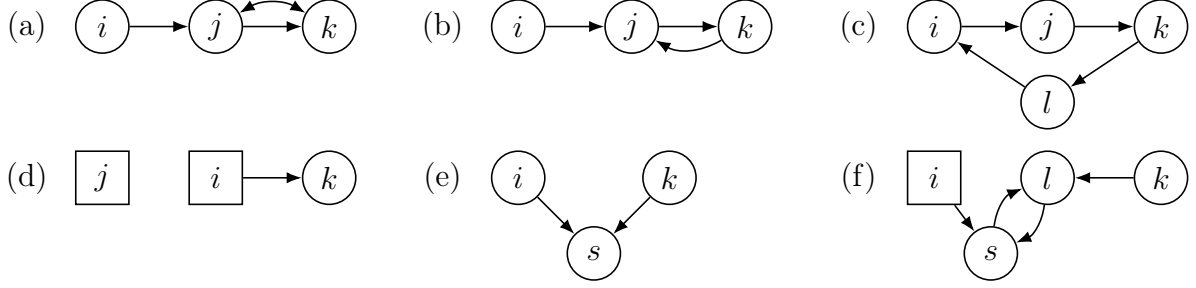


Figure 26: Examples of non-trivial  $\sigma$ -inducing paths in CDMGs. (a) The path  $i \rightarrow j \leftrightarrow k$  is a  $\sigma$ -inducing path between  $i$  and  $k$ . (b) The path  $i \rightarrow j \rightarrow k$  is a  $\sigma$ -inducing path between  $i$  and  $k$ . (c) The path  $i \rightarrow j \rightarrow k$  is a  $\sigma$ -inducing path between  $i$  and  $k$ . The nodes  $i$  and  $k$  cannot be  $\sigma$ -separated by any subset not containing  $i, k$  (indeed,  $i \not\perp^\sigma k$  and  $i \not\perp^\sigma k \mid j$  in all three graphs, and additionally  $i \not\perp^\sigma k \mid l$  and  $i \not\perp^\sigma k \mid \{j, l\}$  in the graph in (c)). (d) The path  $i \rightarrow k$  is  $\sigma$ -inducing, while there is no  $\sigma$ -inducing path between  $i$  and  $j$ , nor between  $j$  and  $k$ . (e) The path  $i \rightarrow s \leftarrow k$  is  $\sigma$ -inducing given  $\{s\}$ . (f) The path  $i \rightarrow s \rightarrow l \leftarrow k$  is  $\sigma$ -inducing given  $\{s\}$ .

**Proposition 11.2.3.** *Let  $G$  be a CDMG with input nodes  $J$ , output nodes  $V^+ = V \dot{\cup} S$  and let  $i \in V$  (but  $i \notin J$ ) and  $j \in V \dot{\cup} J$  be distinct nodes. Then the following are equivalent:*

- (i) *There is a  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  and  $j$ ;*
- (ii) *There is a  $\sigma$ -inducing walk given  $S$  in  $G$  between  $i$  and  $j$ ;*
- (iii)  *$i \not\perp_G^\sigma j \mid S \cup Z$  for all  $Z \subseteq (V \cup J) \setminus \{i, j\}$ ;*
- (iv)  *$i \not\perp_G^\sigma j \mid S \cup Z$  for  $Z = (\text{Anc}_G(\{i, j\} \cup S) \cup J) \setminus \{i, j\}$ .*

*Proof.* The proof is similar to that of Theorem 4.2 in [RS02].

(i)  $\implies$  (ii) is trivial.

(ii)  $\implies$  (iii): Assume the existence of a  $\sigma$ -inducing walk given  $S$  between  $i$  and  $j$  in  $G$ . Let  $Z \subseteq (V \cup J) \setminus \{i, j\}$ . Consider all walks in  $G$  between  $i$  and  $j$  with the property that all colliders on it are in  $\text{Anc}_G(\{i, j\} \cup S \cup Z)$ , and each non-endpoint non-collider on it is not in  $S \cup Z$  or is unblockable. Such walks exist, since the  $\sigma$ -inducing walk is one. Let  $\mu$  be such a walk with a minimal number of colliders. We show that all colliders on  $\mu$  must be in  $\text{Anc}_G(S \cup Z)$ . Suppose on the contrary the existence of a collider  $k$  on  $\mu$  that is not ancestor of  $S \cup Z$ . It is either ancestor of  $i$  or of  $j$ , by assumption. If  $j \in J$ , it cannot be ancestor of  $j$ , and hence must be ancestor of  $i$ . Otherwise, we can assume it to be ancestor of  $i$  without loss of generality. Then there is a directed path  $\pi$  from  $k$  to  $i$  in  $G$  that does not pass through any node of  $S \cup Z$ . Then the subwalk of  $\mu$  between  $k$  and  $j$  can be concatenated with the directed path  $\pi$  into a walk between  $i$  and  $j$  that has the property, but has fewer colliders than  $\mu$ : a contradiction. Therefore,  $\mu$  is  $\sigma$ -open given  $S \cup Z$ . Hence,  $i$  and  $j$  are  $\sigma$ -connected given  $S \cup Z$ .

(iii)  $\implies$  (iv) is trivial.

(iv)  $\implies$  (i): Suppose that  $i$  and  $j$  are  $\sigma$ -connected given  $Z = (\text{Anc}_G(\{i, j\} \cup S) \cup J) \setminus \{i, j\}$ . Let  $\pi$  be a path between  $i$  and  $\{j\} \cup J$  that is  $\sigma$ -open given  $Z$ . The end nodes of  $\pi$  must be  $i$  and  $j$ , because  $J \setminus \{i, j\} \subseteq Z$ . We show that  $\pi$  must be a  $\sigma$ -inducing path given  $S$ . First, all colliders on  $\pi$  are in  $\text{Anc}_G(Z)$ , but not in  $J$ , and hence in  $\text{Anc}_G(\{i, j\} \cup S)$ . Second, let  $k$  be any non-endpoint non-collider on  $\pi$ . Then there must be a directed subpath of  $\pi$  starting at  $k$  that ends either at the first collider on  $\pi$  next to  $k$  or at an end node of  $\pi$ , and hence  $k$  must be in  $Z$ . Since  $\pi$  is  $\sigma$ -open given  $Z$ ,  $k$  must be unblockable. Hence, all non-endpoint non-colliders on  $\pi$  must be unblockable.  $\square$

In words: there is a  $\sigma$ -inducing path between two nodes in a CDMG (provided they are not both input nodes) if and only if the two nodes cannot be  $\sigma$ -separated by any subset of the other nodes.

**Remark 11.2.4.** *Two input nodes cannot be  $\sigma$ -separated by some subset of other nodes. Indeed, the trivial path  $j$  is  $\sigma$ -open as long as we don't condition on  $j$  itself. On the other hand, there may, or may not, be a  $\sigma$ -inducing path given  $S$  between two input nodes. This is why Proposition 11.2.3 does not consider the case  $i, j \in J$ .*

The orientations of the outermost edges on a  $\sigma$ -inducing path contain important information about ancestral relations.

**Lemma 11.2.5.** *Let  $G$  be a CDMG with input nodes  $J$ , output nodes  $V^+ = V \dot{\cup} S$  and let  $i, j \in V \dot{\cup} J$  be distinct. If there exists a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  in  $G$ , and all  $\sigma$ -inducing paths given  $S$  in  $G$  between  $i$  and  $j$  are out of  $j$ , then  $j \in \text{Anc}_G(\{i\} \cup S)$ .*

*Proof.* Let  $\mu$  be a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  in  $G$ . It must be of the form  $i \cdots l \leftarrow j$  (with possibly  $l = i$ ). First we show that  $l$  cannot be in  $\text{Sc}_G(j)$ . If  $l \in \text{Sc}_G(j)$ , then let  $\pi$  be a directed path in  $G$  from  $l$  to  $j$  that is entirely contained in  $\text{Sc}_G(j)$ . Let  $m$  be the node on  $\mu$  closest to  $i$  that is also on  $\pi$  (possibly  $m = l$ ). The subpath of  $\pi$  between  $j$  and  $m$  can be concatenated with the subpath of  $\mu$  between  $m$  and  $i$  into a walk between  $j$  and  $i$ . This must be a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  that is into  $j$  by construction: contradiction. Hence  $l$  cannot be in  $\text{Sc}_G(j)$ .

If  $\mu$  is a directed path all the way to  $i$ , then clearly,  $j \in \text{Anc}_G(\{i\} \cup S)$ . Otherwise, it must contain a collider. Let  $k$  be the collider on  $\mu$  closest to  $j$ .  $k$  must be ancestor of  $i$  or  $j$  or  $S$ . In the first and third cases, clearly  $j \in \text{Anc}_G(\{i\} \cup S)$ . In the second case, all nodes on the subpath of  $\mu$  between  $j$  and  $k$  must be in  $\text{Sc}_G(j)$ , a contradiction.  $\square$

**Lemma 11.2.6.** *Let  $G$  be a CDMG with input nodes  $J$ , output nodes  $V^+ = V \dot{\cup} S$  and let  $i, j \in V \dot{\cup} J$  be distinct. If there exists a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  in  $G$  into  $j$ , and  $i \notin \text{Anc}_G(\{j\} \cup S)$ , then there exists a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  in  $G$  that is both into  $i$  and into  $j$ .*

*Proof.* Let  $\mu$  be a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  in  $G$  into  $j$ . This rules out  $j \in J$ . If  $\mu$  is into  $i$ , we are done. Therefore, suppose it is of the form  $i \rightarrow \cdots * \rightarrow j$ . It

cannot be a directed path, since  $i \notin \text{Anc}_G(\{j\} \cup S)$ . Therefore, there must be a collider  $k$  on  $\mu$  such that  $\mu$  is of the form  $i \rightarrow \dots \rightarrow k \leftarrow \dots \leftarrow j$  (with the subpath between  $i$  and  $k$  directed). Then  $k \in \text{Anc}_G(\{i\})$  (sic!), and hence all nodes on  $\mu$  between  $i$  and  $k$  must be in  $\text{Sc}_G(i)$ . Let  $\pi$  be a directed path in  $G$  from  $k$  to  $i$  that is entirely contained in  $\text{Sc}_G(i)$ . Let  $l$  be the node on  $\mu$  closest to  $j$  that is also on  $\pi$  (possibly  $l = k$ ). Then  $l \neq j$ , because otherwise  $j \in \text{Sc}_G(i)$ , contradicting  $i \notin \text{Anc}_G(\{j\} \cup S)$ . The non-trivial subpath of  $\pi$  between  $i$  and  $l$  can be concatenated with the non-trivial subpath of  $\mu$  between  $l$  and  $j$  into a walk between  $i$  and  $j$ . This must be a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  that is both into  $i$  and into  $j$ .  $\square$

**Lemma 11.2.7.** *Let  $G$  be a CDMG with input nodes  $J$ , output nodes  $V^+ = V \dot{\cup} S$  and let  $i, j \in V \dot{\cup} J$  be distinct. If there is a  $\sigma$ -inducing path in  $G$  given  $S$  between  $i$  and  $j$  that is into  $j$ , and  $k \in \text{Sc}_G(j)$ , then there is a  $\sigma$ -inducing path in  $G$  given  $S$  between  $i$  and  $k$  that is into  $k$ .*

*Proof.* The  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  into  $j$  can be extended by a directed path within  $\text{Sc}_G(k)$  to a  $\sigma$ -inducing walk given  $S$  in  $G$  between  $i$  and  $k$  that is into  $k$ .  $\square$

### 11.3. Partial Ancestral Graphs

It is often convenient when performing causal reasoning to be able to represent a set of CDMGs in a compact way. For this purpose, *partial ancestral graphs (PAGs)* have been introduced [SMR99,Zha06]. In order to deal with possible cycles (in simple SCMs), selection bias and exogenous input nodes, we extend the definition to  $\sigma$ -PAGs.

$\sigma$ -PAGs have two node types, input nodes and output nodes (just like CDMGs). They also have multiple edge types. In addition to the three edge types for CDMGs ( $\rightarrow$ ,  $\leftarrow$ ,  $\leftrightarrow$ ), there is an undirected edge ( $\text{---}$ ) and there are five edge types involving a circle edge mark:  $\leftarrow\circ$ ,  $\circ\leftarrow$ ,  $\circ\rightarrow$ ,  $\rightarrow\circ$ ,  $\circ\text{---}$ . Each edge  $i \ast\ast j$  has two *edge marks*, one at each node, with each edge mark either a *tail*, *arrowhead* or *circle*. For example, the directed edge  $i \rightarrow j$  has a tail at  $i$  and an arrowhead at  $j$ , while the bi-circle edge  $i \circ\text{---} j$  has two circle edge marks. All 9 possible combinations of edge marks can occur on an edge in a  $\sigma$ -PAG. We will make use of the “ $\ast$ ” symbol to denote any of the three edge marks. So the notation  $i \ast\ast j$  can stand for all 9 possible edge types between  $i$  and  $j$ , whereas  $i \leftarrow\ast j$  is shorthand for three possible edge types, as are  $i \text{---}\ast j$  and  $i \circ\ast j$ . Edges of the form  $i \leftarrow\ast j$  and  $j \ast\rightarrow i$  are called *into*  $i$ . Edges of the form  $i \text{---}\ast j$  and  $j \ast\text{---} i$  are called *out of*  $i$ .

In order to define  $\sigma$ -PAGs, we extend the definitions of (directed) walks, (directed) paths and colliders to cover these new edge types.

**Definition 11.3.1.** *Let  $H = (J, V, E)$  be a mixed graph with input nodes  $J$ , output nodes  $V$  and edges  $E$  of the types  $\{\rightarrow, \leftarrow, \leftrightarrow, \leftarrow\circ, \circ\leftarrow, \circ\rightarrow, \rightarrow\circ, \circ\text{---}, \text{---}\}$ .<sup>63</sup> Let  $v, w \in V \cup J$ .*

1. *If there is an edge  $v \ast\ast w$  between  $v$  and  $w$  (of any type), we call  $v$  and  $w$  adjacent in  $H$ .*

<sup>63</sup>Formally, we no longer introduce separate sets to represent the edges of each type, but merge them into the single set  $E$ .

2. Define  $\text{Adj}_H(v)$  to be the nodes adjacent to  $v$  in  $H$ .
3. A triple of distinct nodes  $(a, b, c)$  in  $H$  form a triangle if each pair of nodes in the triple is adjacent in  $H$ .
4. A triple of distinct nodes  $(a, b, c)$  in  $H$  is called unshielded if  $b$  is adjacent to both  $a$  and  $c$  in  $H$ , but  $a$  is not adjacent to  $c$  in  $H$ .
5. A walk between  $v$  and  $w$  in  $H$  is a finite alternating sequence of nodes and edges

$$v = v_0, a_0, v_1, \dots, v_{n-1}, a_{n-1}, v_n = w$$

in  $H$  for some  $n \geq 0$ , i.e. such that for every  $k = 0, \dots, n-1$  we have that  $v_k, v_{k+1} \in V \cup J$  and  $a_k = v_k \text{ ** } v_{k+1} \in E$ , and with end nodes  $v_0 = v$  and  $v_n = w$ .

6. A walk is called a path if no node occurs more than once on the walk.
7. A directed walk (path) from  $v$  to  $w$  in  $H$  is a walk (path) of the form:

$$v = v_0 \longrightarrow v_1 \longrightarrow \dots \longrightarrow v_{n-1} \longrightarrow v_n = w,$$

for some  $n \geq 0$ .

8. A path  $v_0 \text{ ** } \dots \text{ ** } v_n$  in  $H$  is called a possibly directed path from  $v_0$  to  $v_n$  if for each  $i = 1, \dots, n$ , the edge  $v_{i-1} \text{ ** } v_i$  is not into  $v_{i-1}$  and is not out of  $v_i$  (i.e., each edge must be of the form  $v_{i-1} \circ\text{-}\circ v_i$ ,  $v_{i-1} \circ\text{-}\rightarrow v_i$ ,  $v_{i-1} \circ\text{-}\leftarrow v_i$  or  $v_{i-1} \rightarrow v_i$ ).
9. A directed cycle is a directed walk  $v \rightarrow \dots \rightarrow w$ , concatenated with the directed edge  $w \rightarrow v$ .
10. An almost directed cycle is a directed walk  $v \rightarrow \dots \rightarrow w$ , concatenated with the bidirected edge  $w \leftrightarrow v$ .
11. A path of the form  $v_0 \circ\text{-}\circ \dots \circ\text{-}\circ v_n$  (with each edge of the form  $v_i \circ\text{-}\circ v_{i+1}$ ) is called a circle path.
12. A path  $v_0 \text{ ** } \dots \text{ ** } v_n$  in  $H$  is called uncovered if every subsequent triple  $(v_{k-1}, v_k, v_{k+1})$  for  $1 < k < n$  is unshielded.
13. A triple of consecutive nodes  $v_{k-1} \text{ ** } v_k \text{ ** } v_{k+1}$  on a walk in  $H$  is called definite collider if it is of the form  $v_{k-1} \text{ ** } \rightarrow v_k \text{ ** } \leftarrow v_{k+1}$ .
14. A triple of consecutive nodes  $v_{k-1} \text{ ** } v_k \text{ ** } v_{k+1}$  on a walk in  $H$  is called a definite non-collider if it is of the form  $v_{k-1} \text{ ** } v_k \text{ ** } \rightarrow v_{k+1}$  or  $v_{k-1} \text{ ** } \leftarrow v_k \text{ ** } v_{k+1}$ . Furthermore, we also refer to the end nodes  $v_0$  and  $v_n$  of a walk between  $v_0$  and  $v_n$  in  $H$  as definite non-colliders.
15. If there is a directed walk from  $v$  to  $w$  in  $H$  then we say that  $v$  is ancestor of  $w$  in  $H$ , and we write  $v \in \text{Anc}_H(w)$ . For  $W \subseteq V \cup J$ , we define  $\text{Anc}_H(W) := \bigcup_{w \in W} \text{Anc}_H(w)$ .

16. If there is a directed walk from  $v$  to  $w$  in  $H$  then we say that  $w$  is descendant of  $v$  in  $H$ , and we write  $w \in \text{Desc}_H(v)$ . For  $W \subseteq V \cup J$ , we define  $\text{Desc}_H(W) := \bigcup_{w \in W} \text{Desc}_H(w)$ .
17. A walk between  $v, w \in V \cup J$  (with  $v \neq w$ ) in  $H$  is called *definitely inducing* if every non-endpoint node is a definite collider in  $\text{Anc}_H(\{v, w\})$ .
18. A walk between  $v$  and  $w$  in  $H$  is called *definitely open* given  $Z \subseteq V \cup J$  if
  - a) every node on the walk is a definite collider or definite non-collider, and
  - b) every definite collider is in  $\text{Anc}_H(Z)$ , and
  - c) every definite non-collider is not in  $Z$ .

We can now define:<sup>64</sup>

**Definition 11.3.2.** A mixed graph  $H = (J, V, E)$  with input nodes  $J$ , output nodes  $V$  and edges  $E$  of the types  $\{\rightarrow, \leftarrow, \leftarrow\circ, \leftrightarrow, \circ\circ, \circ\rightarrow, \rightarrow\circ, \circ\leftarrow, \text{---}\}$  is called a partial  $\sigma$ -ancestral graph ( $\sigma$ -PAG) if all of the following conditions hold:

1. Between any two distinct nodes there is at most one edge, and there are no edges between a node and itself;
2. The graph contains no directed cycles, no almost directed cycles, and no unshielded triple of the form  $i \ast\rightarrow j \text{---} k$  (“ $\sigma$ -ancestral”);
3. There is no definitely inducing path between any two distinct non-adjacent nodes (“maximal”);
4. No input node is adjacent to any other input node;
5. If there is an edge  $j \ast\ast v$  with  $j \in J, v \in V$  then it must be of the form  $j \text{---} v$ .

$\sigma$ -PAGs are used to represent a set of CDMGs as follows.

**Definition 11.3.3.** Let  $H = (J, V, E)$  be a mixed graph with input nodes  $J$ , output nodes  $V$  and edges  $E$  of the types  $\{\rightarrow, \leftarrow, \leftarrow\circ, \leftrightarrow, \circ\circ, \circ\rightarrow, \rightarrow\circ, \circ\leftarrow, \text{---}\}$ . Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$ . We say that  $H$  represents  $G$  given  $S$  if all of the following hold:

1. Between any two distinct nodes in  $H$  there is at most one edge, and there are no edges between a node and itself;
2. Two nodes  $i, j \in V \cup J$  are adjacent in  $H$  if and only if  $i \neq j$ ,  $\{i, j\} \not\subseteq J$ , and there is a  $\sigma$ -inducing path between  $i$  and  $j$  given  $S$  in  $G$ ;

---

<sup>64</sup>We have incorporated two extensions of the usual definition of PAG [Zha06]: we allow for input nodes, and we have weakened the condition of being ancestral to  $\sigma$ -ancestral. A mixed graph is called *ancestral* if it has no directed cycles, no almost directed cycles, and no triples of the form  $i \ast\rightarrow j \text{---} k$ .



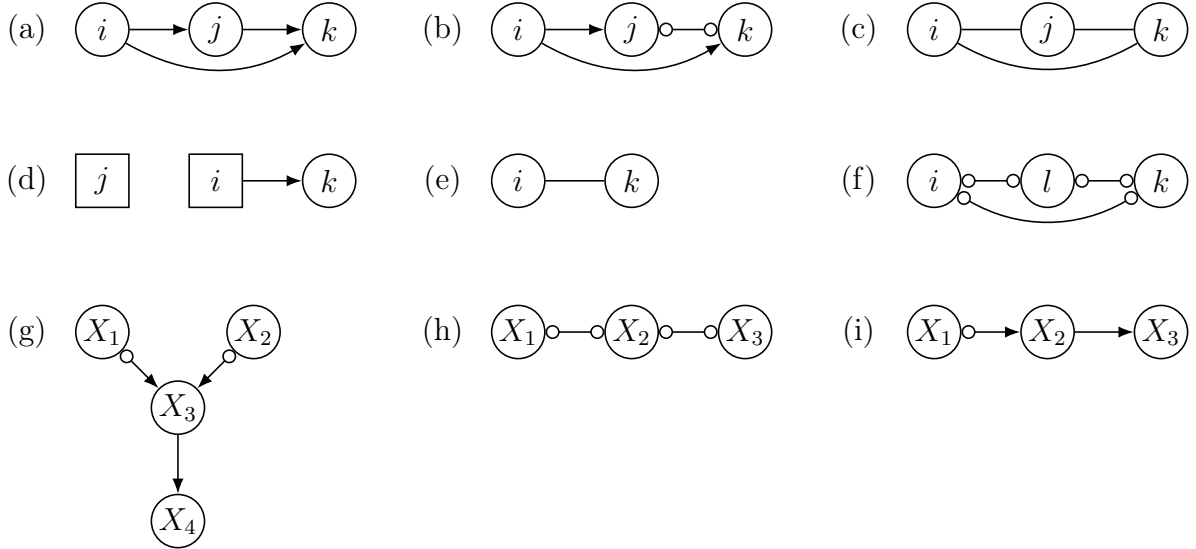


Figure 27: Various example  $\sigma$ -PAGs, representing respectively the CDMGs: (a) in Figure 26(a); (b) in Figures 26(a–b); (c) in Figures 26(a–b); (d) in Figure 26(d); (e) in Figure 26(e); (f) in Figure 26(f); (g) of all Y-structures in Figure 25; (h) of all LCD structures in Figure 23; (i) of all LCD structures in Figure 23.

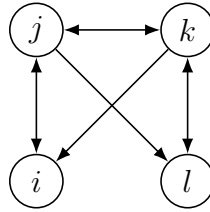


Figure 28: This mixed graph is not a valid  $\sigma$ -PAG, because it is not maximal: it has a definitely inducing path  $i \leftrightarrow j \leftrightarrow k \leftrightarrow l$  while  $i$  and  $l$  are non-adjacent.

3. If  $i \leftarrow^* j$  in  $H$  then  $i \notin \text{Anc}_G(\{j\} \cup S)$ ;
4. If  $i \rightarrow^* j$  in  $H$  then  $i \in \text{Anc}_G(\{j\} \cup S)$ ;
5. If  $j \ast \rightarrow v$  in  $H$  with  $j \in J, v \in V$  then it must be of the form  $j \rightarrow^* v$ .

Hence, adjacencies represent  $\sigma$ -inducing paths, arrowheads represent non-ancestorship (of the adjacent node or  $S$ ), and tails represent ancestorship (of the adjacent node or  $S$ ). Some examples are given in Figure 27. Note in particular that a directed edge in a  $\sigma$ -PAG does not necessarily imply a direct causal relationship (for example, the edge  $i \rightarrow k$  in Figure 27(a)). Figure 28 provides an example of a mixed graph that satisfies all conditions of a  $\sigma$ -PAG except the maximality.

The following property shows that we can no longer unambiguously read off ancestral relations from a  $\sigma$ -PAG that represents a CDMG in case of selection bias ( $S \neq \emptyset$ ).

**Lemma 11.3.4.** *Let  $H$  be a  $\sigma$ -PAG that represents a CDMG  $G$  given  $S$ . For two nodes  $i, j$  in  $H$ :  $i \in \text{Anc}_H(j)$  implies  $i \in \text{Anc}_G(\{j\} \cup S)$ .*

*Proof.* Suppose  $H$  contains a directed path  $i = v_0 \rightarrow \dots \rightarrow v_n = j$ . Note that for all  $k = 0, \dots, n-1$ ,  $v_k \notin \text{Anc}_G(S)$  implies  $v_k \in \text{Anc}_G(v_{k+1})$ . By induction, then,  $v_0 \notin \text{Anc}_G(S)$  implies  $v_0 \in \text{Anc}_G(v_n)$ .  $\square$

The following lemma shows that the orientation of edges in a  $\sigma$ -PAG that represents a CDMG contains information on the orientation of corresponding  $\sigma$ -inducing paths in the CDMG.

**Lemma 11.3.5.** *Let  $H$  be a  $\sigma$ -PAG that represents a CDMG  $G$  given  $S$ . Let  $i, j$  be distinct nodes in  $H$ .*

- (i) *If  $i \star \rightarrow j$  in  $H$ , then there exists a  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  and  $j$  that is into  $j$ .*
- (ii) *If  $i \leftrightarrow j$  in  $H$ , then there exists a  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  and  $j$  that is both into  $i$  and into  $j$ .*

*Proof.* In both cases, there exists a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  in  $G$  because  $i$  and  $j$  are adjacent in  $H$  and  $H$  represents  $G$  given  $S$ . If all  $\sigma$ -inducing paths given  $S$  between  $i$  and  $j$  in  $G$  were out of  $j$ , then by Lemma 11.2.5,  $j \in \text{Anc}_G(\{i\} \cup S)$ , contradicting the orientation  $i \star \rightarrow j$  in  $H$ . Therefore, there must be a  $\sigma$ -inducing path given  $S$  between  $i$  and  $j$  in  $G$  that is into  $j$ . This shows (i). If  $i \leftrightarrow j$  in  $H$ , then application of Lemma 11.2.6 gives (ii).  $\square$

We will frequently use that every mixed graph (of a certain type) that represents a CDMG must be a valid  $\sigma$ -PAG.

**Proposition 11.3.6.** *Let  $H = (J, V, E)$  be a mixed graph with input nodes  $J$ , output nodes  $V$  and edges  $E$  of the types  $\{\rightarrow, \leftarrow, \leftarrow\circ, \leftrightarrow, \circ\circ, \circ\rightarrow, \rightarrow\circ, \circ\leftarrow, \text{---}\}$ . Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \cup S$ . If  $H$  represents  $G$  given  $S$ , then  $H$  is a  $\sigma$ -PAG.*

*Proof.* By assumption: there is at most one edge between two distinct nodes; there are no edges between a node and itself; no input node is adjacent to any other input node; if there is an edge  $j \star \star v$  in  $H$  with  $j \in J, v \in V$  then it must be of the form  $j \rightarrow v$ .

We show that  $H$  is  $\sigma$ -ancestral. First, suppose that  $H$  contained a directed path  $i = v_0 \rightarrow \dots \rightarrow v_n = j$ . By Lemma 11.3.4, then  $i \in \text{Anc}_G(\{j\} \cup S)$ . If  $H$  contained an edge  $j \star \rightarrow i$ , this would imply  $i \notin \text{Anc}_G(\{j\} \cup S)$ , a contradiction. Hence such edges cannot occur. This means that no directed cycles and no almost directed cycles occur in  $H$ . The remaining statement to prove that  $H$  is  $\sigma$ -ancestral follows by Lemma 11.3.7.

We continue to show that  $H$  is maximal. Suppose there is a definitely inducing path  $\mu$  in  $H$  between two distinct nodes  $u, v \in V \cup J$ . We first show that this implies that  $\{u, v\} \not\subseteq J$ . If  $u, v \in J$ , then a definitely inducing path between them cannot consist of a single edge, because all node pairs in  $J$  are non-adjacent in  $H$  by assumption. Hence it must contain at least one non-endpoint node that is not  $u$  or  $v$ . This non-endpoint node must be an ancestor in  $H$  of  $u$  or  $v$ . This is a contradiction, since all edges in  $H$  at

$u$  are out of  $u$ , and all edges in  $H$  at  $v$  are out of  $v$ . Hence, there cannot be a definitely inducing path in  $H$  between two nodes in  $J$ .

Every edge  $i \ast\ast j$  on  $\mu$  corresponds with a  $\sigma$ -inducing path  $\pi_{ij}$  given  $S$  in  $G$  between  $i$  and  $j$ . By Lemma 11.3.5, these  $\sigma$ -inducing paths can be chosen to be into  $i$  if the edge is  $i \leftarrow\ast j$ , into  $j$  if the edge is  $i \ast\rightarrow j$ , and both into  $i$  and  $j$  if the edge is  $i \leftrightarrow j$ . Concatenate all  $\pi_{ij}$  following the edge ordering of  $\mu$  into a walk  $\pi$  in  $G$  between  $u$  and  $v$ . Every non-endpoint node on  $\mu$  is a definite collider on  $\mu$ , by assumption. By construction, these nodes then become colliders on  $\pi$ . Since definite colliders on  $\mu$  are in  $\text{Anc}_H(\{u, v\})$  by assumption, they are in  $\text{Anc}_G(\{u, v\} \cup S)$  by Lemma 11.3.4. All colliders on some  $\pi_{ij}$  are in  $\text{Anc}_G(\{i, j\} \cup S)$ . Hence, they are in  $\text{Anc}_G(\{u, v\} \cup S)$ . So, all colliders on  $\pi$  are in  $\text{Anc}_G(\{u, v\} \cup S)$ . All non-endpoint non-colliders on some  $\pi_{ij}$  are unblockable, and therefore all non-endpoint non-colliders on  $\pi$  are unblockable. Therefore,  $\pi$  is a  $\sigma$ -inducing walk given  $S$  in  $G$ . So there must also be a  $\sigma$ -inducing path given  $S$  in  $G$  between  $u$  and  $v$ . Because  $H$  represents  $G$  given  $S$ , we conclude that  $u, v$  must be adjacent in  $H$ .  $\square$

**Lemma 11.3.7.** *Let  $H = (J, V, E)$  be a mixed graph with input nodes  $J$ , output nodes  $V$  and edges  $E$  of the types  $\{\rightarrow, \leftarrow, \leftarrow\circ, \leftrightarrow, \circ\circ, \circ\rightarrow, \rightarrow\circ, \leftarrow\leftarrow, \rightarrow\rightarrow\}$ . Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$ . If  $H$  represents  $G$  given  $S$ , then no unshielded triple of the form  $i \ast\rightarrow j \text{ --- } k$  can occur in  $H$ .*

*Proof.* By contradiction. Since  $j \in \text{Anc}_G(\{k\} \cup S)$  but  $j \notin \text{Anc}_G(\{i\} \cup S)$ , we must have  $j \in \text{Anc}_G(k)$ . Also,  $k \in \text{Anc}_G(\{j\} \cup S)$ . Since  $j \in \text{Anc}_G(k)$  and  $j \notin \text{Anc}_G(S)$ , we must have  $k \in \text{Anc}_G(j)$ . But then  $j \in \text{Sc}_G(k)$ , which is not possible according to Lemma 11.2.7.  $\square$

The following result shows that every CDMG can be represented by a  $\sigma$ -PAG given some set of selection nodes.

**Proposition 11.3.8.** *Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$ . There exists a  $\sigma$ -PAG, denoted  $\text{PAG}^\sigma(G | S)$ , that represents  $G$  given  $S$ .*

*Proof.* We will construct a mixed graph  $H$  with input nodes  $J$  and output nodes  $V$  as follows. Let two nodes  $i, j \in J \cup V$  be adjacent in  $H$  if and only if  $i \neq j$ ,  $\{i, j\} \not\subseteq J$  and there is a  $\sigma$ -inducing path given  $S$  between  $i, j$  in  $G$ . In that case, orient the edge between  $i$  and  $j$  in  $H$  as follows:

$$\begin{cases} i \text{ --- } j & \text{if } i \in \text{Anc}_G(\{j\} \cup S) \text{ and } j \in \text{Anc}_G(\{i\} \cup S), \\ i \rightarrow j & \text{if } i \in \text{Anc}_G(\{j\} \cup S) \text{ and } j \notin \text{Anc}_G(\{i\} \cup S), \\ i \leftarrow j & \text{if } i \notin \text{Anc}_G(\{j\} \cup S) \text{ and } j \in \text{Anc}_G(\{i\} \cup S), \\ i \leftrightarrow j & \text{if } i \notin \text{Anc}_G(\{j\} \cup S) \text{ and } j \notin \text{Anc}_G(\{i\} \cup S). \end{cases}$$

It is obvious by construction that  $H$  represents  $G$  given  $S$ . It is a valid  $\sigma$ -PAG by Proposition 11.3.6.  $\square$

The  $\sigma$ -PAG constructed in this way is a *maximal  $\sigma$ -ancestral graph* ( $\sigma$ -MAG): it contains no circle edge marks and is therefore maximally informative about ancestral relations (that is, as informative as a  $\sigma$ -PAG can be).

**Proposition 11.3.9.** *Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$  and  $H$  a  $\sigma$ -PAG that represents  $G$  given  $S$ . If  $i \ast \rightarrow j \ast \leftarrow k$  in  $H$  (with  $i, j, k$  distinct nodes) and  $j \in \text{Sc}_G(k)$ , then  $i \ast \leftarrow k$  in  $H$  and  $k \notin \text{Anc}_G(\{i\} \cup S)$ . If, additionally, the edge in  $H$  between  $i$  and  $j$  is  $i \leftrightarrow j$ , then also  $i \notin \text{Anc}_G(\{k\} \cup S)$ .*

*Proof.* Since  $\text{Sc}_G(k) \supseteq \{j, k\}$ ,  $k$  cannot be in  $J$ .

By Lemma 11.3.5,  $i \ast \rightarrow j$  in  $H$  implies the existence of a  $\sigma$ -inducing walk between  $i$  and  $j$  given  $S$  in  $G$  that is into  $j$ . This can be extended by concatenation with a directed path from  $j$  to  $k$  into a  $\sigma$ -inducing walk between  $i$  and  $k$  given  $S$  in  $G$  that is into  $k$ . Hence, there must be an edge  $i \ast \leftarrow k$  in  $H$ .

If  $k \in \text{Anc}_G(\{i\} \cup S)$ , then also  $j \in \text{Anc}_G(\{i\} \cup S)$  because  $j \in \text{Anc}_G(k)$ , a contradiction.

If  $i \leftrightarrow j$  in  $H$ , then  $i \notin \text{Anc}_G(\{j\} \cup S)$ . If  $i \in \text{Anc}_G(\{k\} \cup S)$ , then  $i \in \text{Anc}_G(k)$  because  $i \notin \text{Anc}_G(S)$ ; hence also  $i \in \text{Anc}_G(k) = \text{Anc}_G(j)$ , a contradiction.  $\square$

The following important result states that the existence of a definitely open path in a  $\sigma$ -PAG representing a CDMG may imply the existence of a corresponding  $\sigma$ -open path in the CDMG.

**Proposition 11.3.10.** *Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$  and  $H$  a  $\sigma$ -PAG that represents  $G$  given  $S$ . Let, for  $n \geq 2$ ,*

$$\pi = q_0 \ast \leftarrow q_1 \ast \rightarrow q_2 \leftrightarrow \dots \leftrightarrow q_{n-1} \ast \leftarrow q_n$$

*be a definitely  $Z$ -open path in  $H$ , for  $Z \subseteq (J \cup V) \setminus \{q_0, q_n\}$ , such that  $q_2, \dots, q_{n-1}$  are definite colliders, and  $q_1$  is either a definite non-collider or a definite collider. Then there exists a  $(Z \cup S)$ - $\sigma$ -open path in  $G$  between  $q_0$  and  $q_n$ .*

*Proof.* The nodes  $q_0, q_1, \dots, q_n$  in  $J \cup V$  must all be distinct. For  $i = 1, \dots, n$ , let  $\mu_i$  be a  $\sigma$ -inducing path in  $G$  between  $q_{i-1}$  and  $q_i$  that is into  $q_{i-1}$  if the edge  $q_i \ast \leftarrow q_{i-1}$  on  $\pi$ , and into  $q_i$  if the edge  $q_{i-1} \ast \rightarrow q_i$  on  $\pi$  (see Lemma 11.3.5). These can be concatenated into a walk  $\mu = (\mu_1, \dots, \mu_n)$  in  $G$  between  $q_0$  and  $q_n$ :

$$\overbrace{q_0 \ast \leftarrow \dots \ast \leftarrow q_1}^{\mu_1} \ast \leftarrow \dots \ast \rightarrow q_2 \ast \leftarrow \dots \ast \rightarrow q_{n-2} \ast \leftarrow \dots \ast \rightarrow q_{n-1} \ast \leftarrow \dots \ast \leftarrow q_n^{\mu_n}$$

$$\underbrace{\phantom{q_0 \ast \leftarrow \dots \ast \rightarrow q_2}}_{\mu_2} \quad \underbrace{\phantom{q_{n-2} \ast \leftarrow \dots \ast \rightarrow q_{n-1}}}_{\mu_{n-1}}$$

Let  $Z \subseteq (J \cup V) \setminus \{q_0, q_n\}$ . Since  $\pi$  is assumed to be definitely  $Z$ -open,  $q_2, \dots, q_{n-1}$  are all in  $\text{Anc}_H(Z)$ . We must have that  $q_1 \in \text{Anc}_G(\{q_0, q_2\} \cup Z \cup S)$ , as can be seen by considering the two mutually exclusive cases:

- $q_1$  is a definite collider on  $\pi$ . Then  $q_1 \in \text{Anc}_H(Z)$ . By Lemma 11.3.4,  $q_1 \in \text{Anc}_G(Z \cup S)$ .

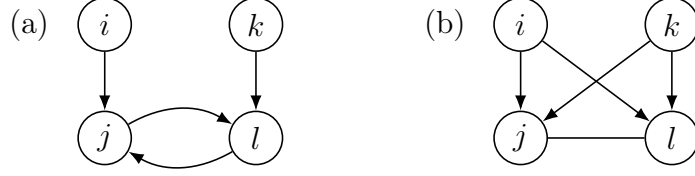


Figure 29: Example that shows that not all definitely open paths in a  $\sigma$ -PAG that represents a CDMG imply the existence of a corresponding  $\sigma$ -open path in the CDMG. (a) CDMG in which  $i \perp^\sigma k$ , represented by the  $\sigma$ -PAG in (b) in which a definitely open path  $i \rightarrow j \leftarrow l \leftarrow k$  exists.

- $q_1$  is a definite non-collider on  $\pi$ . Then  $q_1 \notin Z$ , but  $q_1 \in \text{Anc}_G(\{q_0, q_2\} \cup S)$  because either  $q_0 \leftarrow q_1$  or  $q_1 \rightarrow q_2$  must be on  $\pi$ .

In any case,  $q_1 \in \text{Anc}_G(\{q_0, q_2\} \cup Z \cup S)$ .

Consider all walks in  $G$  between  $q_0$  and  $q_n$  with the property that

1. all colliders on it are in  $\text{Anc}_G(\{q_0, q_n\} \cup Z \cup S)$ , and
2. each non-endpoint non-collider on it is not in  $\{q_0, q_n\} \cup Z \cup S$  or is unblockable.

Such walks exist, since the concatenation  $\mu = (\mu_1, \dots, \mu_n)$  is one, as we will now show.

To show that the first property holds for  $\mu$ , note that  $q_2, \dots, q_{n-1}$  are in  $\text{Anc}_H(Z)$  and hence, by Lemma 11.3.4,  $q_2, \dots, q_{n-1} \in \text{Anc}_G(Z \cup S)$ , a subset of  $\text{Anc}_G(\{q_0, q_n\} \cup Z \cup S)$ . We already saw that  $q_1 \in \text{Anc}_G(\{q_0, q_n\} \cup Z \cup S)$ , which holds in particular if  $q_1$  is a collider on  $\mu$ . Every internal collider on some  $\mu_i$  is in  $\text{Anc}_G(\{q_{i-1}, q_i\} \cup S)$ , and hence also these ‘internal’ colliders on  $\mu$  are in  $\text{Anc}_G(\{q_0, q_n\} \cup Z \cup S)$ .

For the second property, note that all non-endpoint non-colliders on some  $\mu_i$  are unblockable by assumption.  $q_2, \dots, q_{n-1}$  cannot be non-colliders on  $\mu$ . If  $q_1$  is a non-collider on  $\mu$ , then it must be a definite non-collider on  $\pi$  and hence  $q_1 \notin Z$  (and by assumption,  $q_1 \notin S$ ). Hence, the second property also holds for  $\mu$ .

Let  $\nu$  be a walk satisfying both properties above with a minimal number of colliders. We show that all colliders on  $\nu$  must be in  $\text{Anc}_G(Z \cup S)$ .

Suppose on the contrary the existence of a collider  $k$  on  $\nu$  that is not in  $\text{Anc}_G(Z \cup S)$ . It must then be in  $\text{Anc}_G(\{q_0, q_n\})$ , by assumption. Then there exists a directed path in  $G$  from  $k$  to  $q_0$  that does not pass through  $q_n$ , or there exists a directed path from  $k$  to  $q_n$  that does not pass through  $q_0$ . Without loss of generality, assume the former: there is a directed path in  $G$  from  $k$  to  $q_0$  in  $G$  that does not pass through  $Z \cup S \cup \{q_n\}$ . Let  $\pi$  be a shortest path of that type. The directed path  $\pi$  from  $k$  to  $q_0$  can be concatenated with the subwalk of  $\nu$  between  $k$  and  $q_n$  (which is into  $k$ ) into a walk between  $q_0$  and  $q_n$ . This walk has the property, but has fewer colliders than  $\nu$ : a contradiction.

Therefore,  $\nu$  is a  $(Z \cup S)$ - $\sigma$ -open walk in  $G$  between  $q_0$  and  $q_n$ . This means that there must also exist a  $(Z \cup S)$ - $\sigma$ -open path in  $G$  between  $q_0$  and  $q_n$ .  $\square$

## 11.4. Unshielded triples

One of the key steps in the FCI algorithm is the orientation of “unshielded triples”. The following proposition will later be used to “orient” the edges in unshielded triples in a  $\sigma$ -PAG representing a CDMG.

**Proposition 11.4.1.** *Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$  and  $H$  a  $\sigma$ -PAG that represents  $G$  given  $S$ . If  $(i, j, k)$  form an unshielded triple in  $H$  with  $i \in V$  (and  $j, k \in V \cup J$ ), then either*

(i)  $j \notin Z$  for each  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k \mid Z \cup S$ , and  $j \notin \text{Anc}_G(\{i, k\} \cup S)$ ,  
or

(ii)  $j \in Z$  for each  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k \mid Z \cup S$ , and  $j \in \text{Anc}_G(\{i, k\} \cup S)$ .

*Proof.* Case (i):  $j \notin \text{Anc}_G(\{i, k\} \cup S)$ . By orienting the edge marks at  $j$  on the path  $i \ast\ast j \ast\ast k$  in  $H$  into  $i \ast\rightarrow j \ast\leftarrow k$  (if not already oriented that way), we obtain a  $\sigma$ -PAG  $\tilde{H}$  that represents  $G$ . Let  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k \mid Z \cup S$ . If  $j \in Z$ , the path  $i \ast\rightarrow j \ast\leftarrow k$  would be definitely  $Z$ -open in  $\tilde{H}$ . By Proposition 11.3.10, this implies the existence of a  $(Z \cup S)$ - $\sigma$ -open walk in  $G$  between  $i$  and  $k$ , a contradiction. Hence  $j \notin Z$ .

Case (ii):  $j \in \text{Anc}_G(\{i, k\} \cup S)$ . By orienting the edge marks at  $j$  on the path  $i \ast\ast j \ast\ast k$  in  $H$  into  $i \ast\leftarrow j \ast\ast k$  (if  $j \in \text{Anc}_G(\{i\} \cup S)$ ) or  $i \ast\ast j \ast\rightarrow k$  (if  $j \in \text{Anc}_G(\{k\} \cup S)$ ) or  $i \ast\leftarrow j \ast\rightarrow k$  (if  $j \in \text{Anc}_G(\{i\} \cup S) \cap \text{Anc}_G(\{k\} \cup S)$ ), we obtain a  $\sigma$ -PAG  $\tilde{H}$  that represents  $G$ . Let  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k \mid Z \cup S$ . If  $j \notin Z$ , the path  $i \ast\ast j \ast\ast k$  is definitely  $Z$ -open in  $\tilde{H}$ . By Proposition 11.3.10, this implies the existence of a  $(Z \cup S)$ - $\sigma$ -open walk in  $G$  between  $i$  and  $k$ , a contradiction. Hence  $j \in Z$ .  $\square$

Note that in the first case, we can orient the edges as  $i \ast\rightarrow j \ast\leftarrow k$  (if they were not already oriented in this way) to obtain a  $\sigma$ -PAG  $\tilde{H}$  that also represents  $G$ . In the second case, we cannot orient the edges, since we don’t know whether  $j \in \text{Anc}_G(\{i\} \cup S)$  or  $j \in \text{Anc}_G(\{k\} \cup S)$  or both.

## 11.5. Discriminating paths

Another step in the FCI algorithm is related to the notion of “discriminating paths”. This can be considered as an extension of the notion of unshielded triple.

**Definition 11.5.1.** *A path  $\pi = (i, j, q_1, \dots, q_n, k)$  (with  $n \geq 1$ ) in a mixed graph  $H$  is a discriminating path for  $j$  if:*

(i)  $i$  is not adjacent to  $k$  in  $H$ , and

(ii) for  $r = 1, \dots, n$ :  $q_r$  is a definite collider on  $\pi$  and a parent of  $i$  in  $H$ .

Figure 30 illustrates this notion.

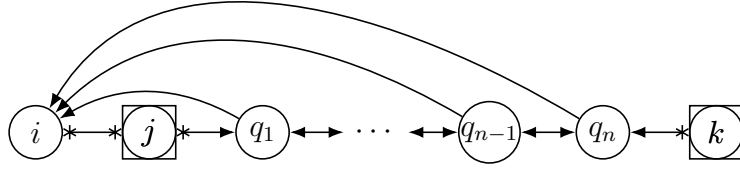


Figure 30: Discriminating path  $(i, j, q_1, \dots, q_n, k)$  for  $j$  between  $i$  and  $k$ . Only  $j$  and  $k$  could be an input node, all other nodes must be output nodes.

**Remark 11.5.2.** *It is instructive to think about a discriminating path rather as a certain collection of paths:*

$$\begin{aligned}
i &\leftarrow q_n \leftarrow^* k \\
i &\leftarrow q_{n-1} \leftrightarrow q_n \leftarrow^* k \\
i &\leftarrow q_{n-2} \leftrightarrow q_{n-1} \leftrightarrow q_n \leftarrow^* k \\
&\vdots \\
i &\leftarrow q_1 \leftrightarrow \dots \leftrightarrow q_{n-1} \leftrightarrow q_n \leftarrow^* k \\
i & \ast \ast j \ast \ast q_1 \leftrightarrow \dots \leftrightarrow q_{n-1} \leftrightarrow q_n \leftarrow^* k
\end{aligned}$$

with the additional requirement that  $i$  and  $k$  are not adjacent.

The following quintessential property of discriminating paths is analogous to that of unshielded triples.

**Proposition 11.5.3.** *Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$  and  $H$  a  $\sigma$ -PAG that represents  $G$  given  $S$ . If  $(i, j, q_1, \dots, q_n, k)$  is a discriminating path in  $H$  for  $j$  between  $i$  and  $k$ , then either*

(i)  $j \notin Z$  for each  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k | Z \cup S$ , and  $j \notin \text{Anc}_G(\{q_1, i\} \cup S)$ ,  
or

(ii)  $j \in Z$  for each  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k | Z \cup S$ , and  $j \in \text{Anc}_G(\{i\} \cup S)$ .

In both cases,  $i \notin \text{Anc}_G(\{j\} \cup S)$ .

*Proof.* Since  $q_1$  is a parent of  $i$  in  $H$ ,  $q_1 \in \text{Anc}_G(\{i\} \cup S)$ . Because  $q_1$  is a definite collider in  $H$ ,  $q_1 \notin \text{Anc}_G(S)$ . Hence  $q_1 \in \text{Anc}_G(i) \setminus \{i\}$ , which means that  $i \in V$ . Also, this implies that  $i \notin \text{Anc}_G(j)$ ; otherwise,  $q_1 \in \text{Anc}_G(j)$  which contradicts  $j \ast \rightarrow q_1$  in  $H$ . Hence  $i \notin \text{Anc}_G(\{j\} \cup S)$ .

We first show that if  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k | Z \cup S$ , then  $\{q_1, \dots, q_n\} \in Z$ . This can be seen from the various subpaths in Remark 11.5.2. From the first path: if  $q_n \notin Z$  then this path would be definitely open in  $H$ , and hence (by Proposition 11.3.10) there must be a  $(Z \cup S)$ - $\sigma$ -open walk in  $G$  between  $i$  and  $k$ , which would contradict  $i \perp_G^\sigma k | Z \cup S$ . Once we have shown that  $\{q_n, q_{n-1}, \dots, q_{n-r+1}\} \in Z$  for  $1 \leq r < n$ , we

see from the  $r + 1$ 'th path that  $q_{n-r} \in Z$  to avoid definitely opening up this path in  $H$ , and thereby avoiding the existence of a  $(Z \cup S)$ - $\sigma$ -open path in  $G$  between  $i$  and  $k$ .

Case (i):  $j \notin \text{Anc}_G(\{q_1, i\} \cup S)$ . By orienting the edge marks at  $j$  on the path  $i \rightsquigarrow j \rightsquigarrow q_1$  in  $H$  into  $i \rightsquigarrow j \leftarrow q_1$  (if not already oriented that way), we obtain a  $\sigma$ -PAG  $\tilde{H}$  that represents  $G$ . Let  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k \mid Z \cup S$ . If  $j \in Z$ , the discriminating path in  $\tilde{H}$  is definitely  $Z$ -open. By Proposition 11.3.10, this implies the existence of a  $(Z \cup S)$ - $\sigma$ -open walk in  $G$  between  $i$  and  $k$ . Contradiction. Hence  $j \notin Z$ .

Case (ii):  $j \in \text{Anc}_G(\{q_1, i\} \cup S)$ . Since  $q_1 \in \text{Anc}_G(i)$ , this implies  $j \in \text{Anc}_G(\{i\} \cup S)$ . We already concluded above that  $i \notin \text{Anc}_G(\{j\} \cup S)$ . By orienting the the edge  $i \rightsquigarrow j$  in  $H$  as  $i \leftarrow j$ , we obtain a  $\sigma$ -PAG  $\tilde{H}$  that represents  $G$ . Let  $Z \subseteq (V \cup J) \setminus \{i, k\}$  such that  $i \perp_G^\sigma k \mid Z \cup S$ . If  $j \notin Z$ , the discriminating path in  $\tilde{H}$  is definitely  $Z$ -open. By Proposition 11.3.10, this implies the existence of a  $(Z \cup S)$ - $\sigma$ -open walk in  $G$  between  $i$  and  $k$ . Contradiction. Hence  $j \in Z$ .  $\square$

Note that in the first case, we can orient  $i \leftrightarrow j \leftrightarrow q_1$  (as far as the edge marks were not already oriented in this way) to obtain a  $\sigma$ -PAG  $\tilde{H}$  that also represents  $G$ . In the second case, we can orient  $i \leftarrow j$  (as far as the edge marks were not already oriented in this way) to obtain a  $\sigma$ -PAG  $\tilde{H}$  that also represents  $G$ .

## 11.6. Independence models and Markov equivalence

We start with an abstract definition of an ‘‘independence model’’, where we extend the common definition to allow for input nodes.

**Definition 11.6.1.** *Given two disjoint sets  $J, V$  (the ‘inputs’ and ‘outputs’, respectively), we call a subset of*

$$\{(A, B, C) : A, B, C \subseteq J \cup V, A \cap J = \emptyset, J \subseteq B \cup C\}$$

*an independence model over  $V \mid J$ . For an element  $(A, B, C)$  of an independence model, we also say that  $C$  separates  $A$  from  $B$ .*

Independence models can be used to encode all (conditional) independences in a Markov kernel.

**Definition 11.6.2.** *For a Markov kernel  $K(X_V \mid X_J) : \mathcal{X}_J \dashrightarrow \mathcal{X}_V$ , we define its independence model to be*

$$\text{IM}(K(X_V \mid X_J)) := \{(A, B, C) : A, B, C \subseteq J \cup V, A \cap J = \emptyset, J \subseteq B \cup C : X_A \perp\!\!\!\perp_{K(X_V \mid X_J)} X_B \mid X_C\},$$

*i.e., the set of all conditional independences (of restricted form) that  $K(X_V \mid X_J)$  satisfies.*

Independence models can also encode all separations in a graph.



**Definition 11.6.3.** For a CDMG  $G$  with input nodes  $J$  and output nodes  $V$ , define its  $\sigma$ -independence model to be

$$\text{IM}_\sigma(G) := \{(A, B, C) : A, B, C \subseteq J \cup V, A \cap J = \emptyset, J \subseteq B \cup C : A \perp_G^\sigma B \mid C\},$$

*i.e.*, the set of all  $\sigma$ -separations (of restricted form) entailed by the graph. Define its  $d$ -independence model to be

$$\text{IM}_d(G) := \{(A, B, C) : A, B, C \subseteq J \cup V, A \cap J = \emptyset, J \subseteq B \cup C : A \perp_G^d B \mid C\},$$

*i.e.*, the set of all  $d$ -separations (of restricted form) entailed by the graph.

Both  $\text{IM}_d(G)$  and  $\text{IM}_\sigma(G)$  are independence models over  $V \mid J$  (with  $J$  the input nodes of  $G$  and  $V$  the output nodes of  $G$ ). For CADMGs,  $\sigma$ -separation is equivalent to  $d$ -separation, and hence, if  $G$  is acyclic, then  $\text{IM}_d(G) = \text{IM}_\sigma(G)$ .

When conditioning on a set of selection variables, we can also define conditional independence models from graphs as follows.

**Definition 11.6.4.** For a CDMG  $G$  with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$ , define its  $\sigma$ -independence model given  $S$  to be

$$\text{IM}_\sigma(G \mid S) := \{(A, B, C) : A, B, C \subseteq J \cup V, A \cap J = \emptyset, J \subseteq B \cup C : A \perp_{G \mid S}^\sigma B \mid C \cup S\},$$

*i.e.*, the set of all  $\sigma$ -separations (of restricted form) involving nodes not in  $S$ , when also conditioning on  $S$ . Define its  $d$ -independence model given  $S$  to be

$$\text{IM}_d(G \mid S) := \{(A, B, C) : A, B, C \subseteq J \cup V, A \cap J = \emptyset, J \subseteq B \cup C : A \perp_{G \mid S}^d B \mid C \cup S\},$$

*i.e.*, the set of all  $d$ -separations (of restricted form) involving nodes not in  $S$ , when also conditioning on  $S$ .

Hence, also  $\text{IM}_\sigma(G \mid S)$  and  $\text{IM}_d(G \mid S)$  are independence models over  $V \mid J$  (with  $J$  the input nodes of  $G$  and  $V$  the output nodes of  $G$  except for those in  $S$ ). For CADMGs,  $\sigma$ -separation is equivalent to  $d$ -separation, and hence, if  $G$  is acyclic, then  $\text{IM}_d(G \mid S) = \text{IM}_\sigma(G \mid S)$ .

The input to the extended FCI algorithm will consist of an independence model over  $V \mid J$ , and its output will consist of a mixed graph with input nodes  $J$  and output nodes  $V$ . One of the nice properties of the extended FCI algorithm is that if the input of FCI is the independence model of a CDMG given  $S$ , then the output of FCI will be a  $\sigma$ -PAG that represents the ‘‘Markov equivalence class’’ of the CDMG (at least if  $S = \emptyset$  or  $J = \emptyset$ ).

**Definition 11.6.5.** Let  $G_1, G_2$  be two CDMGs with input nodes  $J$  and output nodes  $V_1^+ = V \dot{\cup} S_1, V_2^+ = V \dot{\cup} S_2$ , respectively. We call  $G_1$  and  $G_2$   $\sigma$ -Markov equivalent w.r.t.  $V \mid J$  if  $\text{IM}_\sigma(G_1 \mid S_1) = \text{IM}_\sigma(G_2 \mid S_2)$ , and  $d$ -Markov equivalent w.r.t.  $V \mid J$  if  $\text{IM}_d(G_1 \mid S_1) = \text{IM}_d(G_2 \mid S_2)$ .

When exploiting conditional independences for causal discovery, one typically has to make some kind of faithfulness assumption. The faithfulness assumption that we will make for the extended FCI algorithm will be that the (restricted) conditional independences in the conditioned Markov kernel are identical to the (restricted)  $\sigma$ -separations in the observable graph, given  $S$ :

$$\text{IM}(P_M(X_V \mid X_S \in \xi_S, \text{do}(X_J))) = \text{IM}_\sigma(G_{V+|J}(M) \mid S).$$

## 11.7. Skeleton search

The FCI algorithm consists of two phases, the skeleton search phase, which is followed by the orientation phase. In this subsection, we describe the skeleton search phase.

**Definition 11.7.1.** *Given a  $\sigma$ -PAG  $H = (J, V, E)$ , its skeleton is the mixed graph  $\text{skel}(H) := (J, V, F)$  with the same nodes, and with a bicircle edge  $i \circ\!\circ j$  in  $F$  if and only if  $i \ast\ast j$  in  $E$  (i.e., if  $i$  and  $j$  are adjacent in  $H$ ).*

Hence, the only edge type occurring in the skeleton is the bicircle edge. The skeleton has no edge between any pair of input nodes. We will later frequently refer to the unordered pairs of nodes that may be adjacent:

**Definition 11.7.2.** *For input nodes  $J$  and output nodes  $V$ , define*

$$\text{separable}(V|J) := \{\{i, j\} : i \in V, j \in V, i \neq j\} \cup \{(i, j) : i \in V, j \in J\}.$$

The aim of the skeleton search phase of the FCI algorithm is to construct the skeleton of the  $\sigma$ -PAG that represents a CDMG given  $S$  from the  $\sigma$ -independence model given  $S$  of the CDMG. It does this by testing for each separable edge in the skeleton whether it can find any subset of nodes that separates the two nodes. If it finds a separating set between an (unordered) pair of distinct nodes  $\{i, j\}$ , the set is stored as  $\text{sepset}(\{i, j\})$ . The orientation phase later makes use of these separating sets found in the skeleton phase.

A brute-force search over all possible subsets of  $(J \cup V) \setminus \{i, j\}$ , as in Algorithm 1, would be a straightforward solution. However, it is also computationally extremely expensive for all but the smallest cardinalities of  $V$  and  $J$ , and statistically not very reliable in case the separations have to be tested with conditional independence tests on finite data.

To get some inspiration on how to address this, we will first describe the skeleton search phase of the PC algorithm, an ancestor of the FCI algorithm designed for DAGs [SGS00]. The PC skeleton phase (Algorithm 2) searches for separating sets of increasing cardinality. As candidates for a separating set between nodes  $i, j \in H$ , it considers all subsets of  $\text{Adj}_H(i)$  and all subsets of  $\text{Adj}_H(j)$ . The rationale is that if  $G$  is a DAG, then either  $\text{Pa}_G(i) \subseteq \text{Adj}_H(i)$  or  $\text{Pa}_G(j) \subseteq \text{Adj}_H(j)$   $d$ -separates  $i$  from  $j$ .

We can easily extend this to deal with input nodes as well, see Algorithm 3. We may restrict ourselves to search for separating sets that contain all nodes in  $J$  (except  $j$  itself, if  $j \in J$ ). We therefore only need to consider subsets of neighbours of  $i$  that are not in  $J$ , in increasing cardinality.

---

**Algorithm 1** Brute-force skeleton algorithm.

---

1: **Input:** Input node set  $J$ ; output node set  $V$ ; independence model  $I$  over  $V \mid J$   
2: **Output:** mixed graph  $H$  with input nodes  $J$  and output nodes  $V$ ; separating sets  $\text{sepset}$   
3:  $H \leftarrow (J, V, \emptyset)$   
4: **for each**  $(i, j) \in \text{separable}(V \mid J)$  **do**  
5:     add edge  $i \circ\!\!\!\circ j$  to  $H$   
6:     **for each**  $Z \subseteq (V \cup J) \setminus \{i, j\}$  **do**  
7:         **if**  $(i, j, Z) \in I$  **then** ▷ found a separating set  
8:             delete edge  $i \circ\!\!\!\circ j$  from  $H$   
9:              $\text{sepset}(\{i, j\}) \leftarrow Z$   
10:             **break**  
11:         **end if**  
12:     **end for**  
13: **end for**

---

---

**Algorithm 2** Original PC skeleton algorithm.

---

1: **Input:** Output node set  $V$ ; independence model  $I$  over  $V \mid \emptyset$   
2: **Output:** mixed graph  $H$  with output nodes  $V$ ; separating sets  $\text{sepset}$   
3: initialize  $H$  as a complete graph with only bicircle ( $\circ\!\!\!\circ$ ) edges  
4:  $n \leftarrow 0$   
5: **repeat**  
6:     **repeat**  
7:         select  $i, j \in V$  with  $i \circ\!\!\!\circ j$  in  $H$  and  $\#(\text{Adj}_H(i) \setminus \{j\}) \geq n$   
8:         select a subset  $Z \subseteq \text{Adj}_H(i) \setminus \{j\}$  of cardinality  $n$   
9:         **if**  $(i, j, Z) \in I$  **then**  
10:             delete edge  $i \circ\!\!\!\circ j$  from  $H$   
11:              $\text{sepset}(\{i, j\}) \leftarrow Z$   
12:         **end if**  
13:     **until** no more such tuples  $(i, j, Z)$  can be selected  
14:      $n \leftarrow n + 1$   
15: **until** for all  $i, j \in V$  with  $i \circ\!\!\!\circ j$  in  $H$ ,  $\#(\text{Adj}_H(i) \setminus \{j\}) < n$

---

---

**Algorithm 3** Extended PC skeleton algorithm  $\text{PCskeleton}(J, V, I)$ .

---

```

1: Input: Input node set  $J$ ; output node set  $V$ ; independence model  $I$  over  $V \mid J$ 
2: Output: mixed graph  $H$  with input nodes  $J$  and output nodes  $V$ ; separating sets
   sepset
3:  $H \leftarrow (J, V, \emptyset)$ 
4: for each  $(i, j) \in \text{separable}(V \mid J)$  do
5:   add edge  $i \circ\!\!\!\circ j$  to  $H$ 
6: end for
7:  $n \leftarrow 0$ 
8: repeat
9:   repeat
10:    select  $i \in V, j \in V \cup J$  with  $i \circ\!\!\!\circ j$  in  $H$  and  $\#(\text{Adj}_H(i) \setminus (J \cup \{j\})) \geq n$ 
11:    select a subset  $Z \subseteq \text{Adj}_H(i) \setminus (J \cup \{j\})$  of cardinality  $n$ 
12:    if  $(i, j, Z \cup (J \setminus \{j\})) \in I$  then
13:      delete edge  $i \circ\!\!\!\circ j$  from  $H$ 
14:       $\text{sepset}(\{i, j\}) \leftarrow Z \cup (J \setminus \{j\})$ 
15:    end if
16:   until no more such tuples  $(i, j, Z)$  can be selected
17:    $n \leftarrow n + 1$ 
18: until for all  $i \in V, j \in V \cup J$  with  $i \circ\!\!\!\circ j$  in  $H$ ,  $\#(\text{Adj}_H(i) \setminus (J \cup \{j\})) < n$ 

```

---

The underlying idea may no longer hold if  $G$  is a CADMG or a CDMG, or in case of selection bias. The original proposal (which motivated the somewhat optimistic adjective “Fast” in the name of the FCI algorithm [SMR95]) replaces the subsets of adjacent nodes by so-called “Possible-D-Sep” sets.<sup>65</sup> Before we can define these, we need some definitions and theory.

**Definition 11.7.3.** Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$ . For distinct nodes  $i, j \in V \dot{\cup} J$ , define  $\text{SEP}_G(i, j)$  as the set of nodes  $k \in V$  such that  $k \neq i$  and there is a walk  $\pi$  in  $G$  between  $i$  and  $k$  such that every node on  $\pi$  is in  $\text{Anc}_G(\{i, j\} \cup S)$ , and every non-endpoint non-collider on  $\pi$  is unblockable.

The name  $\text{SEP}_G(i, j)$  is motivated by the following property.

**Proposition 11.7.4.** Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$ . Let  $i \in V$  and  $j \in J \cup V$  be distinct nodes. If there exists no  $\sigma$ -inducing walk given  $S$  between  $i, j$  in  $G$ , then  $\text{SEP}_G(i, j) \cap \{i, j\} = \emptyset$  and  $i \perp_G^\sigma j \mid \text{SEP}_G(i, j) \cup (J \setminus \{j\}) \cup S$ .

*Proof.* By definition,  $\text{SEP}_G(i, j) \subseteq \text{Anc}_G(\{i, j\} \cup S)$  and  $i \notin \text{SEP}_G(i, j)$ . If  $j \in \text{SEP}_G(i, j)$ , the walk  $\pi$  in Definition 11.7.3 between  $i$  and  $j$  must be  $\sigma$ -inducing given  $S$ . Since there exists no  $\sigma$ -inducing walk given  $S$  between  $i, j$  in  $G$  by assumption,  $j \notin \text{SEP}_G(i, j)$ .

---

<sup>65</sup>An alternative search strategy was proposed that can be considerably faster in practice, for which it can be shown that the corresponding FCI+ algorithm is of polynomial-time complexity in the number of nodes, as long as the degree of the DPAG is bounded [CMH13].

We prove the  $\sigma$ -separation by contradiction. Write  $Z := \text{SEP}_G(i, j) \cup (J \setminus \{j\}) \cup S$ . Suppose there exists a path in  $G$  between  $i$  and  $\{j\} \cup J$  that is  $\sigma$ -open given  $Z$ . It cannot end in a node in  $J \setminus \{j\}$ , and hence it must end in  $j$ . Let  $\pi = v_0 \ast\ast \dots \ast\ast v_n$  with  $v_0 = i$  and  $v_n \in j$  be such a path consisting of a minimal number of nodes. Every node on  $\pi$  must be in  $\text{Anc}_G(\{i, j\} \cup S) \cup J$ , since by Lemma 9.10.2, each node on  $\pi$  is in  $\text{Anc}_G((\{i, j\} \setminus J) \cup Z)$ , and  $Z \subseteq \text{Anc}_G(\{i, j\} \cup S) \cup J$ .  $\pi$  can only contain nodes in  $J$  as endnodes (otherwise we could shorten the path). In other words, the only node on  $\pi$  that could be in  $J$  is  $j$ .

Denote the subwalk of  $\pi$  from  $v_a$  to (and including)  $v_b$  by  $\pi(v_a, v_b)$ , for  $0 \leq a \leq b \leq n$ . We will show that for all  $k = 1, \dots, n$ ,  $\pi(v_0, v_k)$  has the property that every non-endpoint non-collider on it is unblockable. The property trivially holds for  $k = 1$ . Suppose it holds for  $k < n$ . Since  $v_0, \dots, v_k$  are all in  $\text{Anc}_G(\{i, j\} \cup S)$ , and all non-endpoint non-colliders on  $\pi(v_0, v_k)$  are unblockable, we conclude that  $v_k \in \text{SEP}_G(i, j)$ . If  $v_k$  is a non-collider on  $\pi(v_0, v_{k+1})$ , it must be unblockable, because  $\pi(v_0, v_{k+1})$  is  $Z$ - $\sigma$ -open and  $v_k \in \text{SEP}_G(i, j) \subseteq Z$ . So the property also holds for  $k + 1$ .

In particular, we can conclude that  $j = v_n \in \text{SEP}_G(i, j)$ , a contradiction.  $\square$

In practice, one does not know the set  $\text{SEP}_G(i, j)$  if  $G$  is unknown. One can, however, easily obtain a ‘bound’ on this set by identifying a superset.

**Definition 11.7.5.** Let  $H_0 = (J, V, E)$  be a mixed graph with input nodes  $J$ , output nodes  $V$  and edges  $E$  of the types  $\{\rightarrow, \leftarrow, \leftarrow\circ, \leftrightarrow, \circ\circ, \circ\rightarrow, \rightarrow\circ, \text{---}\}$ . For  $i \in V, j \in J \cup V$  distinct, we define  $\text{posSEP}_{H_0}(i, j) \subseteq V$  to consist of those nodes  $k \in V$  such that  $k \notin \{i, j\}$  and there is a path between  $i$  and  $k$  in  $H_0$  such that for every subsequent triple  $a \ast\ast b \ast\ast c$  on the path, either the triple is a definite collider in  $H_0$ , or a triangle in  $H_0$ .

We can now show that:

**Lemma 11.7.6.** Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \cup S$ . Let  $H_0$  be the mixed graph constructed by lines 3–8 of Algorithm 4 when run on the independence model  $\text{IM}_\sigma(G|S)$ . Then, for  $i \in V$  and  $j \in J \cup V$  distinct nodes, if there is no  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  and  $j$ , then  $\text{SEP}_G(i, j) \subseteq \text{posSEP}_{H_0}(i, j)$ .

*Proof.* Note that the skeleton of  $H_0$  is a supergraph of the skeleton of  $\text{PAG}^\sigma(G|S)$ , i.e., every adjacency in  $\text{PAG}^\sigma(G|S)$  must also be an adjacency in  $H_0$ . Let  $k \in \text{SEP}_G(i, j)$ . Since there is no  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  and  $j$ ,  $\{i, j\} \cap \text{SEP}_G(i, j) = \emptyset$  by Proposition 11.7.4, and hence  $k \neq i$  and  $k \neq j$ . By definition, there is a path  $\pi$  in  $G$  between  $i$  and  $k$  with the property that every node on  $\pi$  is in  $\text{Anc}_G(\{i, j\} \cup S) \cap V$ , and every non-endpoint non-collider on  $\pi$  is unblockable. There is a corresponding path  $\pi'$  in  $H_0$  between  $i$  and  $k$  that consists of the same sequence of nodes (but may have different edges between the nodes).

Consider a subsequent triple  $a \ast\ast b \ast\ast c$  on  $\pi$ . Suppose  $b$  is a collider on  $\pi$ . In  $H_0$ ,  $a$  must also be adjacent to  $b$ , and  $b$  to  $c$ . If  $a$  and  $c$  are not adjacent in  $H_0$ , then a separating set for  $a$  and  $c$  was found during the construction of  $H_0$ , or  $a$  and  $c$  are both in  $J$ . The latter would be a contradiction. Therefore, it will have been oriented

as a definite collider in  $H_0$  according to Proposition 11.4.1. Otherwise (if  $a$  and  $c$  are adjacent in  $H_0$ ), it forms a triangle in  $H_0$ . If  $b$  is a non-collider on  $\pi$ , then it must be unblockable. But in that case,  $a \ast\ast b \ast\ast c$  is a  $\sigma$ -inducing walk given  $S$  between  $a$  and  $c$  in  $G$ . Hence,  $(a, b, c)$  will form a triangle in  $H_0$ . Therefore,  $k \in \text{posSEP}_{H_0}(i, j)$ .  $\square$

---

**Algorithm 4** Extended FCI skeleton algorithm  $\text{FCIskeleton}(J, V, I)$ .

---

```

1: Input: Input node set  $J$ ; output node set  $V$ ; independence model  $I$  over  $V \mid J$ 
2: Output: mixed graph  $H$  with input nodes  $J$  and output nodes  $V$ ; separating sets
   sepset
3:  $(H_0, \text{sepset}) \leftarrow \text{PCskeleton}(J, V, I)$ 
4: for each unshielded triple  $(i, j, k)$  in  $H_0$  with  $i, j \in V$  do ▷ orient colliders
5:   if  $j \notin \text{sepset}(\{i, k\})$  then
6:     orient it as  $i \ast\rightarrow j \leftarrow\ast k$ 
7:   end if
8: end for
9: for each  $i \in V, j \in V \cup J$  with  $i \neq j$  do ▷ construct Possible-D-Sep sets
10:  calculate  $\text{posSEP}_{H_0}(i, j)$ 
11: end for
12:  $H \leftarrow (J, V, \{i \circ\text{--} \circ j : i \ast\ast j \in H_0\})$  ▷ forget orientations
13:  $n \leftarrow 0$ 
14: repeat ▷ search for separating sets
15:   repeat
16:     select  $i \in V, j \in V \cup J$  with  $i \circ\text{--} \circ j$  in  $H$  and  $\#(\text{posSEP}_{H_0}(i, j)) \geq n$ 
17:     select a subset  $Z \subseteq \text{posSEP}_{H_0}(i, j)$  of cardinality  $n$ 
18:     if  $(i, j, Z \cup (J \setminus \{j\})) \in I$  then
19:       delete edge  $i \circ\text{--} \circ j$  from  $H$ 
20:        $\text{sepset}(\{i, j\}) \leftarrow Z \cup (J \setminus \{j\})$ 
21:     end if
22:   until no more such tuples  $(i, j, Z)$  can be selected
23:    $n \leftarrow n + 1$ 
24: until for all  $i \in V, j \in V \cup J$  with  $i \circ\text{--} \circ j$  in  $H$ ,  $\#(\text{posSEP}_{H_0}(i, j)) < n$ 

```

---

So we do not need to search over all possible subsets of  $(J \cup V) \setminus \{i, j\}$  for a separating set between  $i, j$ , but only the subsets in  $\text{posSEP}_{H_0}(i, j)$ . The skeleton phase of the extended FCI algorithm is described in Algorithm 4. It is an extension of the original FCI skeleton search [SMR95] to deal with input nodes. It first runs the extended PC skeleton phase (Algorithm 3) and orients unshielded triples, obtaining a directed mixed graph  $H_0$ . Then it calculates the sets  $\text{posSEP}_{H_0}(i, j)$  for all distinct pairs  $(i, j)$  with  $i \in V, j \in J \cup V$ . If  $i \perp_G^\sigma j \mid Z \cup S$  for some  $Z \subseteq (J \cup V) \setminus \{i, j\}$ , then  $i \perp_G^\sigma j \mid Z^* \cup S$  for some  $Z^* \subseteq \text{posSEP}_{H_0}(i, j)$ . Since some of the oriented colliders may be incorrect in  $H_0$  (because the PC skeleton phase may not have found all separating sets), it removes all orientations from  $H_0$  and then continues with a more extensive search for separating

sets, similar to how the PC skeleton phase is done, but now using  $\text{posSEP}_{H_0}(i, j)$  instead of  $\text{Adj}_{H_0}(i) \setminus \{j\}$  to find candidate nodes for the separating set.

**Theorem 11.7.7.** *The extended FCI skeleton algorithm (Algorithm 4) is sound: if its input consists of the  $\sigma$ -independence model  $\text{IM}_\sigma(G|S)$  given  $S$  of a CDMG  $G$ , then its output will be  $\text{skel}(\text{PAG}^\sigma(G|S))$ . Furthermore,  $i \perp^\sigma j \mid \text{sepset}(\{i, j\})$  for all  $i \in V, j \in V \cup J$  for node pairs  $(i, j) \in \text{separable}(V|J)$  that are non-adjacent in the skeleton.*

*Proof.* Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \dot{\cup} S$  and  $I = \text{IM}_\sigma(G|S)$  its  $\sigma$ -independence model given  $S$ . By Lemma 11.7.6, the mixed graph  $H_0$  constructed by lines 3–8 has the property that for distinct  $i \in V, j \in J \cup V$ , if there is no  $\sigma$ -inducing path given  $S$  in  $G$  between  $i$  and  $j$ , then  $\text{SEP}_G(i, j) \subseteq \text{posSEP}_{H_0}(i, j)$ . Thus, if  $i \in V$  and  $j \in J \cup V$  are not adjacent in  $\text{skel}(\text{PAG}^\sigma(G|S))$ , then we are guaranteed that for some subset  $Z \subseteq \text{posSEP}_{H_0}(i, j)$ , we have that  $i \perp_G^\sigma j \mid (Z \cup J) \setminus \{j\} \cup S$ . The graph  $H$  constructed in lines 12–24 will be a mixed graph with input nodes  $J$  and output nodes  $V$  that only has bicircle edges. It has an edge between any pair of distinct nodes  $i, j \in J \cup V$  if and only if  $\{i, j\} \not\subseteq J$  and there is no set  $Z \subseteq (J \cup V) \setminus \{i, j\}$  such that  $i \perp_G^\sigma j \mid Z \cup S$ . Furthermore, if there is no edge in  $H$  between a pair of distinct nodes  $i, j \in J \cup V$ , then either  $\{i, j\} \subseteq J$ , or  $i \perp_G^\sigma j \mid \text{sepset}(\{i, j\}) \cup S$ . By Proposition 11.2.3, this implies that two distinct nodes  $i, j \in J \cup V$  are adjacent in  $H$  at this stage if and only if there is a  $\sigma$ -inducing walk given  $S$  in  $G$  between  $i$  and  $j$ . Hence  $H$  must be  $\text{skel}(\text{PAG}^\sigma(G|S))$ .  $\square$

## 11.8. FCI Algorithm

We are now ready to describe a causal inference algorithm that is an extension of the original Fast Causal Inference (FCI) algorithm of [SMR95] to deal with input nodes and cycles. It is presented as Algorithm 5. Its input is an independence model over  $V|J$ , where  $V$  and  $J$  are index sets of output and input nodes, respectively. Its output is a mixed graph with input nodes  $J$  and output nodes  $V$ . It starts with a *skeleton phase* (line 3, see Algorithm 4) that is aimed at deducing the adjacencies between the nodes, and to find sets that separate two separable node pairs. Then, it runs various *orientation rules* (lines 4–16) that iteratively orient circle edge marks into tails and arrowheads. Note that by convention, the labeled nodes within each orientation rule are assumed to be distinct (for example, in  $\mathcal{R}0$ , it is implicitly assumed that  $i \neq j \neq k \neq i$ ). For the special case  $J = \emptyset$ , the algorithm reduces to the standard formulation of the FCI algorithm [Zha08].<sup>66</sup>

<sup>66</sup>Compared to the standard formulation of [Zha08], which assumes no input nodes, we have adapted the skeleton search phase (by starting with a mixed graph that contains no edges between input nodes, limiting the paths in the calculation of the  $\text{posSEP}_{H_0}(i, j)$  sets to output nodes only, and by including all input nodes, except  $j$  itself, into the separating set). Furthermore, we added step 4 to orient the edges between input and output nodes. The formulation of orientation rules  $\mathcal{R}0, \mathcal{R}1, \mathcal{R}3, \mathcal{R}7$  is slightly different for this extended version (as these would not be valid in case both  $i, k \in J$ ), and the rest of the algorithm is unchanged.

---

**Algorithm 5** Extended FCI Algorithm.

---

1: **Input:** Input node set  $J$ ; output node set  $V$ ; independence model  $I$  over  $V \mid J$   
2: **Output:** mixed graph  $H$  with input nodes  $J$  and output nodes  $V$   
3:  $(H, \text{sepset}) \leftarrow \text{FCISkeleton}(J, V, I)$   
4: **for each** edge  $j \circ\!\!\circ v$  in  $H$  with  $j \in J, v \in V$  **do**  
5:     orient  $j \circ\!\!\circ v$   
6: **end for**  
7: **repeat**  
    $\mathcal{R}0$  if  $i \ast\ast j \ast\ast k$  in  $H$  with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  
       $i \ast\rightarrow j \ast\leftarrow k$  if  $j \notin \text{sepset}(\{i, k\})$   
8: **until** this orientation rule is not applicable  
9: **repeat**  
    $\mathcal{R}1$  if  $i \ast\!\circ j \ast\leftarrow k$  in  $H$  with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  
       $i \leftarrow j$   
    $\mathcal{R}2$  if  $i \rightarrow j \ast\rightarrow k$  or  $i \ast\rightarrow j \rightarrow k$  in  $H$ , and  $i \ast\!\circ k$  in  $H$ , then orient  $i \ast\rightarrow k$   
    $\mathcal{R}3$  if  $i \ast\rightarrow j \ast\leftarrow k$  and  $i \ast\!\circ l \circ\ast k$  and  $l \ast\!\circ j$  in  $H$  with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  $l \ast\rightarrow j$   
    $\mathcal{R}4$  if  $(i, j, q_1, \dots, q_n, k)$  is a discriminating path in  $H$  for  $j$ , and if  $i \ast\!\circ j$  in  $H$ , then orient  $i \leftarrow j$  if  $j \in \text{sepset}(\{i, k\})$  and orient  $i \leftrightarrow j \leftrightarrow q_1$  if  $j \notin \text{sepset}(\{i, k\})$   
10: **until** none of these orientation rules is applicable  
11: **repeat**  
    $\mathcal{R}5$  if  $i \circ\!\!\circ j$  in  $H$ , and there is an uncovered circle path  $i \circ\!\!\circ k \circ\!\!\circ \dots \circ\!\!\circ l \circ\!\!\circ j$  in  $H$  such that  $i$  is not adjacent to  $l$  and  $j$  is not adjacent to  $k$ , then orient  
       $i \text{ --- } k \text{ --- } \dots \text{ --- } l \text{ --- } j \text{ --- } i$   
12: **until** this orientation rule is not applicable  
13: **repeat**  
    $\mathcal{R}6$  if  $i \text{ --- } j \circ\ast k$  in  $H$ , then orient  $j \text{ --- } \ast k$   
    $\mathcal{R}7$  if  $i \ast\!\circ j \circ\text{---} k$  in  $H$  with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  
       $i \ast\text{---} j$   
14: **until** none of these orientation rules is applicable  
15: **repeat**  
    $\mathcal{R}8$  if  $i \rightarrow j \rightarrow k$  in  $H$ , and  $i \circ\rightarrow k$  in  $H$ , then orient  $i \rightarrow k$   
    $\mathcal{R}9$  if  $i \circ\rightarrow k$ , and  $\pi = (i, j, \dots, k)$  is an uncovered possibly directed path in  $H$  from  $i$  to  $k$  such that  $j$  and  $k$  are not adjacent in  $H$ , then orient  $i \rightarrow k$   
    $\mathcal{R}10$  if  $i \circ\rightarrow k$  in  $H$ ,  $j \rightarrow k \leftarrow l$  in  $H$ ,  $\pi_1$  is a uncovered possibly directed path in  $H$  from  $i$  to  $j$ , and  $\pi_2$  is a uncovered possibly directed path in  $H$  from  $i$  to  $l$ , then let  $u_1$  be the node adjacent to  $i$  on  $\pi_1$  (possibly  $u_1 = j$ ) and  $u_2$  the node adjacent to  $i$  on  $\pi_2$  (possible  $u_2 = l$ ); if  $u_1 \neq u_2$ , and  $u_1$  and  $u_2$  are not adjacent in  $H$ , then orient  $i \rightarrow k$   
16: **until** none of these orientation rules is applicable

---



**Theorem 11.8.1** (Extended FCI Soundness). *The Extended FCI algorithm (Algorithm 5) is sound: if its input consists of the  $\sigma$ -independence model  $\text{IM}_\sigma(G|S)$  of a CDMG  $G$  given  $S$ , then its output will be a valid  $\sigma$ -PAG  $H$  that represents  $G$  given  $S$ .*

*Proof.* Let  $G$  be a CDMG with input nodes  $J$  and output nodes  $V^+ = V \cup S$  and  $I = \text{IM}_\sigma(G|S)$  its  $\sigma$ -independence model given  $S$ .

The skeleton phase in line 3, which invokes Algorithm 4, is sound (Theorem 11.7.7). That is, it computes  $H = \text{skel}(\text{PAG}^\sigma(G|S))$ , and `sepset` will contain a separating set of  $i$  and  $j$  for every edge  $i \circ\!\!\circ j$  absent in  $H$  with  $i \in V, j \in J \cup V, i \neq j$ . The next step, line 4, orients the edges between input and output nodes in  $H$ . The result is then a valid  $\sigma$ -PAG that represents  $G$  given  $S$ . The rest of the proof proceeds by induction.

We show for each of the orientation rules that under the assumption that the current  $H$  is a valid  $\sigma$ -PAG that represents  $G$  given  $S$ , applying the rule yields an updated  $H$  that is still a valid  $\sigma$ -PAG that represents  $G$  given  $S$ . Some of the rules ( $\mathcal{R}1, \mathcal{R}3, \mathcal{R}5, \mathcal{R}6, \mathcal{R}7, \mathcal{R}9, \mathcal{R}10$ ) assume that rule  $\mathcal{R}0$  has been exhaustively applied, which is the reason that rule  $\mathcal{R}0$  is performed before the other orientation rules are performed. Additionally, rule  $\mathcal{R}6$  assumes that rule  $\mathcal{R}5$  has been exhaustively applied.

In the following, we will always assume that the antecedent of the rule holds for a mixed graph  $H$  that is a valid  $\sigma$ -PAG that represents CDMG  $G$  given  $S$ . This implies, in particular, that if  $v \leftarrow^* w$  in  $H$  or  $v \circ\!\!\circ^* w$  in  $H$ , then  $v \notin J$ .

$\mathcal{R}0$  “If  $i \leftarrow^* j \leftarrow^* k$  in  $H$ , with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  $i \leftarrow^* j \leftarrow^* k$  if  $j \notin \text{sepset}(\{i, k\})$ .”

It follows from Proposition 11.4.1 that  $H$  still represents  $G$  given  $S$  after the orientation of the unshielded collider.

$\mathcal{R}1$  “If  $i \circ\!\!\circ j \leftarrow^* k$  in  $H$  with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  $i \leftarrow j$ .”

We first show that  $j \in \text{Anc}_G(i)$ . Since the triple  $(i, j, k)$  is an unshielded triple in  $H$  with  $i \notin J$ , but has not been oriented as a collider by  $\mathcal{R}0$ , we conclude that  $j \in \text{sepset}(\{i, k\})$ . By Proposition 11.4.1,  $j \in \text{Anc}_G(\{i, k\} \cup S)$ . Since  $H$  represents  $G$  and  $j \leftarrow^* k$  in  $H$ ,  $j \notin \text{Anc}_G(\{k\} \cup S)$ . Therefore,  $j \in \text{Anc}_G(i)$ . Thus if we orient  $i \leftarrow j$ , the mixed graph still represents  $G$  given  $S$ . In particular, the orientation  $i \leftarrow j \leftarrow^* k$  with  $i, k$  non-adjacent cannot occur. Therefore, if we orient  $i \leftarrow j$ , the resulting mixed graph  $H$  will still represent  $G$ .

$\mathcal{R}2$  “If  $i \rightarrow j \leftarrow^* k$  or  $i \leftarrow^* j \rightarrow k$  in  $H$ , and  $i \circ\!\!\circ k$  in  $H$ , then orient  $i \leftarrow^* k$ .”

By the antecedent of the rule, and since  $H$  represents  $G$ , we have  $k \notin \text{Anc}_G(\{j\} \cup S)$ . In case  $i \rightarrow j \leftarrow^* k$ , we have  $i \in \text{Anc}_G(\{j\} \cup S)$ , so if  $k$  were in  $\text{Anc}_G(i)$ , it would follow that  $k \in \text{Anc}_G(\{j\} \cup S)$ , a contradiction. In case  $i \leftarrow^* j \rightarrow k$ , we have  $j \in \text{Anc}_G(\{k\} \cup S)$ , so if  $k$  were in  $\text{Anc}_G(i)$ , it would follow that  $j \in \text{Anc}_G(\{i\} \cup S)$ , a contradiction. Hence, in both cases, we must have  $k \notin \text{Anc}_G(i)$ . In both cases, we also have  $k \notin \text{Anc}_G(S)$  because of the arrowhead on  $j \leftarrow^* k$ . Therefore, after orienting  $i \leftarrow^* k$  in  $H$ , the resulting mixed graph still represents  $G$ .

- $\mathcal{R}3$  “If  $i \ast \rightarrow j \leftarrow \ast k$  and  $i \ast \circ l \circ \ast k$  and  $l \ast \circ j$  in  $H$  with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  $l \ast \rightarrow j$ .”  
 Since  $(i, l, k)$  is an unshielded triple in  $H$  with  $i \notin J$  that was not oriented as a collider by  $\mathcal{R}0$ , we must have that  $l \in \text{Anc}_G(\{i, k\} \cup S)$  by Proposition 11.4.1. Assume, for the sake of contradiction, that  $j \in \text{Anc}_G(\{l\} \cup S)$ . Then  $j \in \text{Anc}_G(\{i, k\} \cup S)$ . This contradicts that  $j \notin \text{Anc}_G(\{i, k\} \cup S)$  from  $i \ast \rightarrow j \leftarrow \ast k$  in  $H$ . Hence,  $j \notin \text{Anc}_G(l \cup S)$ . After orienting  $l \ast \rightarrow j$  in  $H$ , the resulting mixed graph still represents  $G$ .
- $\mathcal{R}4$  “If  $(i, j, q_1, \dots, q_n, k)$  is a discriminating path in  $H$  for  $j$ , and if  $i \ast \circ j$  in  $H$ , then orient  $i \leftarrow j$  if  $j \in \text{sepset}(\{i, k\})$  and orient  $i \leftrightarrow j \leftrightarrow q_1$  if  $j \notin \text{sepset}(\{i, k\})$ .”  
 It follows immediately from Proposition 11.5.3 that applying this rule yields an updated mixed graph that still represents  $G$ .
- $\mathcal{R}5$  “If  $i \circ \circ j$  in  $H$ , and there is an uncovered circle path  $i \circ \circ k \circ \circ \dots \circ \circ l \circ \circ j$  in  $H$  such that  $i$  is not adjacent to  $l$  and  $j$  is not adjacent to  $k$ , then orient  $i \text{ --- } k \text{ --- } \dots \text{ --- } l \text{ --- } j \text{ --- } i$ .”  
 There is an uncovered cycle consisting of (at least 4)  $\circ \circ$  edges in  $H$ . Lemma 11.8.2 implies that each node on the uncovered cycle must be in  $\text{Anc}_G(S)$ . Hence we can orient all edges on the cycle as undirected, and this yields a mixed graph that still represents  $G$  given  $S$ .
- $\mathcal{R}6$  “If  $i \text{ --- } j \circ \ast k$  in  $H$ , then orient  $j \text{ --- } \ast k$ .”  
 Only  $\mathcal{R}5$  could have introduced the undirected edge. In that case, we know that both  $i$  and  $j$  are in  $\text{Anc}_G(S)$  from Lemma 11.8.2. Hence we can orient  $j \text{ --- } \ast k$  and the updated mixed graph will still represent  $G$  given  $S$ .
- $\mathcal{R}7$  “If  $i \ast \circ j \circ \text{ --- } k$  in  $H$  with  $i \notin J$ , and  $i$  and  $k$  are not adjacent in  $H$ , then orient  $i \ast \text{ --- } j$ .”  
 Suppose after orienting  $i \ast \rightarrow j \circ \text{ --- } k$ , the mixed graph would still represent  $G$  given  $S$ . If  $j \in \text{Anc}_G(\{k\} \cup S)$ , then we could further orient  $i \ast \rightarrow j \text{ --- } k$  and the mixed graph would still represent  $G$  given  $S$ , yielding a contradiction because an unshielded triple of that form cannot occur. So  $j \notin \text{Anc}_G(\{k\} \cup S)$ , and we can orient  $i \ast \rightarrow j \leftarrow \ast k$  to obtain a mixed graph that still represents  $G$  given  $S$ . But then we have an unshielded collider  $(i, j, k)$  with  $i \notin J$  that should have been oriented by  $\mathcal{R}0$ , another contradiction. Hence we can orient  $i \ast \text{ --- } j$  and the resulting mixed graph will still represent  $G$  given  $S$ .
- $\mathcal{R}8$  “If  $i \rightarrow j \rightarrow k$  in  $H$ , and  $i \circ \rightarrow k$  in  $H$ , then orient  $i \rightarrow k$ .”  
 It follows immediately from Lemma 11.3.4 that  $i \in \text{Anc}_G(\{k\} \cup S)$ . Thus, applying this rule yields an updated mixed graph that still represents  $G$ .
- $\mathcal{R}9$  “If  $i \circ \rightarrow k$ , and  $\pi = (i, j, \dots, k)$  is an uncovered possibly directed path in  $H$  from  $i$  to  $k$  such that  $j$  and  $k$  are not adjacent in  $H$ , then orient  $i \rightarrow k$ .”  
 Note that  $j \notin J$  as either  $i \ast \rightarrow j$  or  $i \ast \circ j$ . The unshielded triple  $(j, i, k)$ , with  $j \notin J$ , was not oriented as a collider by  $\mathcal{R}0$ , and therefore  $i \in \text{Anc}_G(\{j, k\} \cup S)$ .

If  $i \in \text{Anc}_G(j)$ , then Lemma 11.8.3 then states that  $i \in \text{Anc}_G(k)$ . Therefore, we conclude that  $i \in \text{Anc}_G(\{k\} \cup S)$ . Applying the rule therefore yields an updated mixed graph that still represents  $G$ .

**R10** “Suppose that  $i \circ \rightarrow k$  in  $H$ ,  $j \rightarrow k \leftarrow l$  in  $H$ ,  $\pi_1$  is a uncovered possibly directed path in  $H$  from  $i$  to  $j$ , and  $\pi_2$  is an uncovered possibly directed path in  $H$  from  $i$  to  $l$ . Let  $u_1$  be the node adjacent to  $i$  on  $\pi_1$  (possibly  $u_1 = j$ ) and  $u_2$  the node adjacent to  $i$  on  $\pi_2$  (possibly  $u_2 = l$ ). If  $u_1 \neq u_2$ , and  $u_1$  and  $u_2$  are not adjacent in  $H$ , then orient  $i \rightarrow k$ .”

Because  $u_1$  lies on a possibly directed path starting at  $i$ , it cannot be in  $J$ ; the same holds for  $u_2$ . The unshielded triple  $(u_1, i, u_2)$ , with  $u_1, u_2 \notin J$ , was not oriented by rule  $\mathcal{R}0$ , which implies that  $i \in \text{Anc}_G(\{u_1, u_2\} \cup S)$ . If  $i \in \text{Anc}_G(u_1)$ , Lemma 11.8.3 gives that  $i \in \text{Anc}_G(j)$ . If  $i \in \text{Anc}_G(u_2)$ , Lemma 11.8.3 gives that  $i \in \text{Anc}_G(l)$ . In both cases,  $i \in \text{Anc}_G(\{k\} \cup S)$ , because  $j \rightarrow k \leftarrow l$  in  $H$ . So we conclude that  $i \in \text{Anc}_G(\{k\} \cup S)$  and can orient  $i \rightarrow k$  to obtain an updated mixed graph that still represents  $G$  given  $S$ .

Since each of these orientation rules leaves the skeleton (adjacencies) of  $H$  invariant, and after each orientation rule,  $H$  still represents  $G$ ,  $H$  remains a valid  $\sigma$ -PAG throughout the orientation phase by Proposition 11.3.6.  $\square$

The following two lemmata are applicable after the initial phase of the extended FCI algorithm, once rule  $\mathcal{R}0$  has been exhaustively applied. The first concerns uncovered circle paths, the second uncovered possibly directed paths.

**Lemma 11.8.2.** *Let  $H$  be a mixed graph that represents  $G$  given  $S$  in which rule  $\mathcal{R}0$  has been exhaustively applied. If  $i \circ \circ j$  in  $H$ , and there is an uncovered circle path  $i \circ \circ k \circ \circ \dots \circ \circ l \circ \circ j$  in  $H$  such that  $i$  is not adjacent to  $l$  and  $j$  is not adjacent to  $k$ , then every node on the path is in  $\text{Anc}_G(S)$ .*

*Proof.* There is an uncovered cycle consisting of (at least 4)  $\circ \circ$  edges in  $H$ . Note that none of the nodes on the cycle can be in  $J$ , as each node has two or more circle edge marks. Suppose we orient one of the circles into an arrowhead and the new mixed graph still represents  $G$  given  $S$ . Then we could make use of rule  $\mathcal{R}1$  repeatedly to orient the whole cycle as a directed cycle, and the resulting mixed graph should still represent  $G$  given  $S$ . However, that would be a contradiction, since it contains a directed cycle. Hence, any circle edge mark on the cycle that we orient as an arrowhead would yield a mixed graph that no longer represents  $G$  given  $S$ . In other words, each node on the cycle must be ancestor in  $G$  of its neighboring node, or of the selection set  $S$ . This implies that every node on the cycle must be in  $\text{Anc}_G(S)$ . Indeed, if a given node on the cycle is not in  $\text{Anc}_G(S)$ , then it must be ancestor in  $G$  of its neighbors. Its neighbors then cannot be in  $\text{Anc}_G(S)$  either, and therefore must be ancestor in  $G$  of their neighbors. But then we would have an unshielded triple of nodes that are in the same strongly connected component of  $G$ : a contradiction with Lemma 11.2.2.  $\square$

**Lemma 11.8.3.** *Let  $H$  be a mixed graph that represents  $G$  given  $S$  in which rule  $\mathcal{R}0$  has been exhaustively applied. Let  $v_1 \rightsquigarrow \dots \rightsquigarrow v_n$  be an uncovered possibly directed path from  $v_1$  to  $v_n$  (with  $n \geq 2$ ) in  $H$ . If there is an edge  $k \rightsquigarrow v_1$  in  $H$  and if  $v_1 \in \text{Anc}_G(\{v_2\} \cup S)$ , then we conclude that  $v_1 \in \text{Anc}_G(v_n)$ .*

*Proof.* Note that  $v_1, v_2, \dots, v_n$  are all not in  $J$ , because each of them must have an arrow head or circle edge mark. We can orient  $v_1 \leftarrow v_2$  because of the assumption that  $v_1 \in \text{Anc}_G(\{v_2\} \cup S)$ . Note that  $v_1 \in \text{Anc}_G(v_2)$  because  $v_1 \notin \text{Anc}_G(S)$  due to the arrowhead at  $v_1$  on the edge  $k \rightsquigarrow v_1$ . If  $n = 2$ , then we immediately conclude that  $v_1 \in \text{Anc}_G(v_n)$ . So assume  $n > 2$ . We distinguish two cases:  $v_2 \in \text{Anc}_G(v_1)$  and  $v_2 \notin \text{Anc}_G(v_1)$ .

If  $v_2 \notin \text{Anc}_G(v_1)$ , we can orient  $v_1 \rightarrow v_2$ , since  $v_2 \in \text{Anc}_G(S)$  would imply  $v_1 \in \text{Anc}_G(S)$ , a contradiction. By the same reasoning as in the proof of rule  $\mathcal{R}1$ , we conclude that  $v_2 \in \text{Anc}_G(v_3)$ , and that we can orient  $v_2 \rightarrow v_3$ . We can iterate this reasoning subsequently on all remaining edges on the uncovered possibly directed path to deduce that  $v_i \in \text{Anc}_G(v_{i+1})$  and we can orient  $v_i \rightarrow v_{i+1}$ , for  $i = 3, \dots, n-1$ . This leads to the conclusion that  $v_1 \in \text{Anc}_G(v_n)$ .

If  $v_2 \in \text{Anc}_G(v_1)$ , that is  $v_2 \in \text{Sc}_G(v_1)$ , then we can orient  $v_1 \leftarrow v_2$ . By Lemma 11.2.7, this implies that  $k$  and  $v_2$  must be adjacent in  $H$  as well. This edge can be oriented as  $k \rightsquigarrow v_2$ , because  $v_2 \in \text{Anc}_G(\{k\} \cup S)$  would imply  $v_1 \in \text{Anc}_G(\{k\} \cup S)$ , a contradiction. The assumption  $v_2 \notin \text{Anc}_G(\{v_3\} \cup S)$  gives a contradiction: we could then orient  $v_2 \leftarrow v_3$ , obtaining an unshielded triple  $v_1 \leftarrow v_2 \leftarrow v_3$ . Hence  $v_2 \in \text{Anc}_G(\{v_3\} \cup S)$ , and we can orient  $v_2 \rightarrow v_3$ . Because of the arrowhead in  $k \rightsquigarrow v_2$  at  $v_2$ ,  $v_2 \notin \text{Anc}_G(S)$ , and hence  $v_2 \in \text{Anc}_G(v_3)$ . If  $v_3 \in \text{Anc}_G(\{v_2\} \cup S)$ , then  $v_3 \in \text{Sc}_G(v_2)$  (as  $v_3 \in \text{Anc}_G(S) \setminus \text{Anc}_G(v_2)$  would contradict  $v_2 \notin \text{Anc}_G(S)$ ). But then  $v_1$  and  $v_3$  must lie in the same strongly connected component of  $G$ , and should therefore be adjacent in  $H$  by Lemma 11.2.2, which they are not. Hence  $v_3 \notin \text{Anc}_G(\{v_2\} \cup S)$  and we can orient  $v_2 \rightarrow v_3$ . The reasoning now proceeds as in the previous case, and leads to the conclusion that  $v_2 \in \text{Anc}_G(v_n)$ . Because  $v_1 \in \text{Anc}_G(v_2)$ , also  $v_1 \in \text{Anc}_G(v_n)$ .  $\square$

The soundness of the Extended FCI algorithm immediately implies its consistency when using consistent conditional independence tests.

**Corollary 11.8.4** (Extended FCI Consistency). *Let  $M = (J, V^+, W, \mathcal{X}, P, f)$  be a simple SCM with endogenous variables  $V^+ = V \dot{\cup} S \dot{\cup} L$  (note that we allow additional latent endogenous variables  $L$  here). Let  $\xi_S \subseteq \mathcal{X}_S$  be a measurable set with  $\mathbb{P}_M(X_S \in \xi_S \mid \text{do}(X_J = x_J)) > 0$  for all  $x_j \in \mathcal{X}_J$ . Assume that we have access to infinitely many samples distributed according to the marginal Markov kernel of  $M$  on  $V$  after selecting on  $S$ :*

$$P_M(X_V \mid X_S \in \xi_S, \text{do}(X_J = x_J)).$$

*Assume also that the following faithfulness assumption holds:*

$$\text{IM}(P_M(X_V \mid X_S \in \xi_S, \text{do}(X_J = x_J))) = \text{IM}_\sigma(G_{V \cup S \mid J}(M) \mid S).$$

When using asymptotically consistent conditional independence tests (on the i.i.d. samples of the Markov kernel  $P_M(X_V | X_S \in \xi_S, \text{do}(X_J = x_J))$ ), the Extended FCI algorithm (Algorithm 5) provides an asymptotically consistent estimate  $\hat{H}$  of the  $\sigma$ -PAG  $\text{FCI}(\text{IM}_\sigma(G_{VUS|J}(M) | S))$ , which represents  $G_{VUS|J}(M)$  given  $S$ .

*Proof.* The asymptotic consistency of the conditional independence tests means that the probability of a wrong test result (either Type I or Type II error) vanishes asymptotically. Since  $\text{IM}(P_M(X_V | X_S \in \xi_S, \text{do}(X_J = x_J)))$  consists of finitely many conditional independence statements, the test results will agree completely with this conditional independence model with arbitrarily high probability given sufficiently many samples. By the faithfulness assumption, this conditional independence model agrees with  $\text{IM}_\sigma(G_{VUS|J}(M) | S)$ . By the soundness of the Extended FCI algorithm (Theorem 11.8.1), the  $\sigma$ -PAG that FCI outputs will equal  $\text{FCI}(\text{IM}_\sigma(G_{VUS|J}(M) | S))$  with arbitrarily high probability given sufficiently many samples.  $\square$

Because conditional independence tests are not uniformly consistent (as there is no upper bound on the number of samples needed to distinguish an arbitrarily weak dependence from an independence, without additional assumptions), also the Extended FCI algorithm is not uniformly consistent. In other words, it is not known in advance how many samples will be needed to yield a reliable result.

## 11.9. Completeness

For acyclic  $G$  without input nodes (that is, if  $G$  is an ADMG), the FCI algorithm was shown to be complete [Zha08] in the sense that all edge marks that could possibly be oriented based on the information in  $\text{IM}_\sigma(G | S)$  will be oriented. Using results on the characterization of Markov equivalence classes of maximal ancestral graphs [ARSZ05], it can additionally be shown that the  $\sigma$ -PAG output by FCI represents the Markov equivalence class of  $G$  with respect to  $V$  in case  $G$  is an ADMG. By employing acyclifications, these results have been extended to cyclic  $G$  [MC20], still without input nodes, and with the additional assumption of no selection bias ( $S = \emptyset$ ). The known completeness results have very long proofs, and we will therefore not provide these here. Instead, we will only formulate these results, and refer the interested reader to the original papers for the proofs.

We make the following assumption in order to state the known completeness results.<sup>67</sup>

**Assumption 11.9.1.** *Given node set  $V$ , let:*

- $\mathcal{G}_V$  be the set of all pairs  $(G, S)$  of acyclic DMGs  $G$  with output nodes  $V^+ = V \dot{\cup} S$  for some disjoint set  $S$  ('acyclicity'), or

---

<sup>67</sup>While we believe that it is straightforward to extend the completeness results to allow for both cycles and selection bias, it is known that the Extended FCI algorithm (Algorithm 5) is incomplete when also allowing for input nodes.

- $\mathcal{G}_V$  be the set of all pairs  $(G, \emptyset)$  of DMGs  $G$  with output nodes  $V^+ = V \dot{\cup} \emptyset$  ('no selection bias').

In both cases, all DMGs in  $\mathcal{G}_V$  have  $J = \emptyset$  ('no input nodes').

In the absence of input nodes (for  $J = \emptyset$ ), we can interpret FCI as a mapping that maps an independence model (on  $V$ ) to a mixed graph (with nodes  $V$ ). By Theorem 11.8.1, it maps  $\text{IM}_\sigma(G|S)$ , the  $\sigma$ -independence model of a DMG  $G$  given  $S$ , to a  $\sigma$ -PAG  $\text{FCI}(\text{IM}_\sigma(G|S))$  that represents  $G$  given  $S$ . Additionally, the following (incomplete) completeness results are known.

**Theorem 11.9.2** (Some FCI completeness results). *Under Assumption 11.9.1, the Extended FCI algorithm (Algorithm 5) is:*

- (i) arrowhead complete: for all  $(G, S) \in \mathcal{G}_V$ , for all  $i \neq j \in V$ : there is an arrowhead  $i \leftarrow^* j$  in  $\text{FCI}(\text{IM}_\sigma(G|S))$  if  $i \notin \text{Anc}_{\tilde{G}}(\{j\} \cup S)$  for all  $(\tilde{G}, \tilde{S}) \in \mathcal{G}_V$  with  $\text{IM}_\sigma(\tilde{G}|\tilde{S}) = \text{IM}_\sigma(G|S)$ ;
- (ii) tail complete: for all  $(G, S) \in \mathcal{G}_V$ , for all  $i \neq j \in V$ : there is a tail  $i \rightarrow^* j$  in  $\text{FCI}(\text{IM}_\sigma(G|S))$  if  $i \in \text{Anc}_{\tilde{G}}(\{j\} \cup S)$  for all  $(\tilde{G}, \tilde{S}) \in \mathcal{G}_V$  with  $\text{IM}_\sigma(\tilde{G}|\tilde{S}) = \text{IM}_\sigma(G|S)$ ;
- (iii) Markov complete: for all  $(G_1, S_1) \in \mathcal{G}_V$  and  $(G_2, S_2) \in \mathcal{G}_V$ :  $\text{IM}_\sigma(G_1|S_1) = \text{IM}_\sigma(G_2|S_2)$  if and only if  $\text{FCI}(\text{IM}_\sigma(G_1|S_1)) = \text{FCI}(\text{IM}_\sigma(G_2|S_2))$ .

*Proof.* The first two claims are proved in [Zha08] under the additional assumption of acyclicity. In [MC20] it is explained how the characterization of the Markov equivalence classes of [ARSZ05] can be used to then prove the third claim under that additional assumption. Furthermore, in [MC20] it is shown how to generalize these results to the cyclic setting by employing acyclifications, but only under the additional assumption of no selection bias.  $\square$

Arrowhead and tail completeness express that the  $\sigma$ -PAG output by FCI is maximally oriented: any arrowhead or tail that could possibly be deduced from  $\text{IM}_\sigma(G|S)$ , will have been oriented as such in the  $\sigma$ -PAG. The soundness and Markov completeness properties together imply that the  $\sigma$ -PAG output by FCI, when given as input the  $\sigma$ -independence model of a directed mixed graph given some set of latent selection nodes, represents the  $\sigma$ -Markov equivalence class of  $G$  with respect to the observed nodes. In other words, FCI provides a *graphical characterization* of the  $\sigma$ -Markov equivalence class.

## A. Appendix: Measure Theoretic Probability

This appendix provides a crash course (or rather, refresher) of concepts from measure theoretic probability.

### A.1. Why Measure Theory?

#### Discrete and absolute continuous distributions are not general enough

**Example A.1.1** (Simple example of a non-discrete non-absolute-continuous distribution). Consider a uniformly distributed random variable on the interval  $\mathcal{X} := [0, 1]$ , i.e.  $X \sim \mathcal{U}[0, 1]$ , which has probability density:

$$p(x) = \mathbb{1}_{[0,1]}(x).$$

Consider an exact copy of  $X$ , which we call  $Y := X$ , on  $\mathcal{Y} := [0, 1]$ . Now consider the joint distribution of  $(X, Y)$  on  $\mathcal{X} \times \mathcal{Y} = [0, 1]^2$ . Then only values on the diagonal  $\Delta := \{(x, x) \mid x \in [0, 1]\}$  can be realized by  $(X, Y)$ . This simple distribution on  $[0, 1]^2$  is not discrete (as it can attain uncountably many values), and it is also not absolute continuous, since we have:  $\int_{\Delta} dx dy = 0$ , i.e. the (2-dimensional) area of the (1-dimensional) line is zero. This implies that any density function  $p$  would satisfy:  $\int_{\Delta} p(x, y) dx dy = 0$  as well. This is in contrast to the fact that a probability distribution should always be normalized:

$$1 = P((X, Y) \in \Delta) = \int_{\Delta} p(x, y) dx dy.$$

Note that we don't need a probability density to be able to assign probabilities to subsets  $D \subseteq [0, 1]^2$ . We can just use the push-forward map:

$$(X, Y) : [0, 1] \rightarrow [0, 1] \times [0, 1], \quad x \mapsto (x, x).$$

and compute:

$$P((X, Y) \in D) = P(\{x \in [0, 1] \mid (x, x) \in D\}),$$

where  $P$  on the right here denotes the uniform distribution on  $[0, 1]$ .

**Notation A.1.2** (Unifying the notations to measure theoretic ones). Let  $X$  be a random variable taking values in space  $\mathcal{X}$  and with probability distribution  $P$ . Let  $F : \mathcal{X} \rightarrow \mathbb{R}$  be a function. Then we will change the notations for expectation values as follows.

1. Let  $X$  be a discrete random variable with probability mass function  $p$ . Then define:

$$\begin{aligned} \mathbb{E}[F(X)] &= \sum_{x \in \mathcal{X}} F(x) \cdot p(x) \\ &=: \int F(x) P(dx) \\ &=: \int F(x) dP(x) \\ &=: \int F dP. \end{aligned}$$

We will consider sums to be special cases of measure integrals.

2. Let  $X$  be a absolute continuous random variable with probability density function  $p$ . Then define:

$$\begin{aligned}\mathbb{E}[F(X)] &= \int_{\mathcal{X}} F(x) \cdot p(x) dx \\ &=: \int F(x) P(dx) \\ &=: \int F(x) dP(x) \\ &=: \int F dP.\end{aligned}$$

So both cases can be unified with the 3 commonly used notations:

$$\mathbb{E}[F(X)] = \int F dP = \int F(x) dP(x) = \int F(x) P(dx).$$

Note that in both cases we also can write:  $P(A) = \int \mathbb{1}_A dP$ .

**Exercise A.1.3.** Show that the following relation holds:

$$\int F(x) P(dx) = \int z P^F(dz).$$

### Defining probability distributions on all subsets is too general

**Remark A.1.4.** When we want to work with a (probability) measure  $\mu$  we at least want to require that it is countably additive, i.e. that for pairwise disjoint subsets  $A_n \subseteq \mathcal{X}$ ,  $n \in \mathbb{N}$ , we have:

$$\mu \left( \bigcup_{n \in \mathbb{N}} A_n \right) = \sum_{n \in \mathbb{N}} \mu(A_n).$$

We will see below that if we do not restrict the subsets  $A_n$  in some way we will encounter strange behaviour.

**Theorem A.1.5** (Vitali, non-existence of Lebesgue measure on all subsets). *There does NOT exist a measure  $\lambda$  on  $[0, 1]$  such that:*

1.  $\lambda$  can measure every subset  $A \subseteq [0, 1]$ , and:
2.  $\lambda([a, b]) = b - a$  for all  $a \leq b$  with  $a, b \in [0, 1]$ .

*In other words, there does NOT exist a uniform distribution on  $[0, 1]$  that can consistently assign values to all subsets.*

*But: such a measure with property 2. exists on the set  $\mathcal{B}_{[0,1]}$  of all so called Borel subsets (or even Lebesgue subsets) of  $[0, 1]$ . Similar statements hold for higher dimensions and  $\mathbb{R}^D$  and higher dimensional volumes.*



**Example A.1.6** (Vitali set). Consider the following equivalence relation on  $[0, 1]$ :

$$r_1 \sim r_2 \quad : \iff \quad r_2 - r_1 \in \mathbb{Q}.$$

Let  $[0, 1]/\sim$  be the set of equivalence classes. By the axiom of choice there exists a representative system  $V \subseteq [0, 1]$  for  $[0, 1]/\sim$ . This means that the map:

$$V \rightarrow [0, 1]/\sim, \quad v \mapsto [v],$$

is bijective.  $V$  is called Vitali set and we claim that  $V$  is not Lebesgue-measurable. For this let  $q \in \mathbb{Q}$  and consider the subset:

$$V_q := V + q := \{v + q \mid v \in V\} \subseteq \mathbb{R}.$$

Let  $\mathcal{Q} := [-1, 1] \cap \mathbb{Q}$ . Note that  $[0, 1]$  and  $V_q$  are uncountable while  $\mathbb{Q}$  and  $\mathcal{Q}$  are countably infinite. We then have the inclusions:

$$[0, 1] \subseteq \bigcup_{q \in \mathcal{Q}} V_q \subseteq [-1, 2].$$

The right inclusion is clear as:

$$V + [-1, 1] \subseteq [0, 1] + [-1, 1] \subseteq [-1, 2].$$

For the left inclusion let  $x \in [0, 1]$ . By construction there exists a  $v \in V$  such that  $v \sim x$ . So  $x - v \in \mathbb{Q}$ . Since  $x, v \in [0, 1]$  we also have that  $x - v \in [-1, 1]$ . So  $q := x - v \in [-1, 1] \cap \mathbb{Q} = \mathcal{Q}$ . This shows that  $x \in V_q$  for a  $q \in \mathcal{Q}$ . Thus both inclusions are shown.

If we now assumed that  $V$  would be Lebesgue-measurable then every  $V_q$  would be as well as a translated version of  $V$ . We then would get that:  $\lambda(V_q) = \lambda(V)$  for every  $q \in \mathcal{Q}$ . So we would get:

$$1 = \lambda([0, 1]) \leq \lambda\left(\bigcup_{q \in \mathcal{Q}} V_q\right) \leq \lambda([-1, 2]) = 3,$$

which implies:

$$[1, 3] \ni \lambda\left(\bigcup_{q \in \mathcal{Q}} V_q\right) = \sum_{q \in \mathcal{Q}} \lambda(V_q) = \sum_{q \in \mathcal{Q}} \lambda(V),$$

which is contradictory. Indeed,  $\lambda(V) = 0$  can be ruled out as the sum would sum up to  $0 \notin [1, 3]$ . But also  $\lambda(V) > 0$  can be ruled out as this would sum up to  $\infty \notin [1, 3]$ . So the Vitali set  $V$  can not be Lebesgue-measurable.

**Theorem A.1.7** (Banach-Tarski paradox). The 3-dimensional unit ball  $B_1(z) = \{x \in \mathbb{R}^3 \mid \|x - z\| \leq 1\}$  centered at  $z \in \mathbb{R}^3$  can be partitioned into a finite number of disjoint sets  $A_1, \dots, A_K$  (e.g.  $K = 5$ ) such that each can then be rotated and translated in  $\mathbb{R}^3$  such that they form TWO 3-dimensional unit balls  $B_1(y_1)$  and  $B_1(y_2)$ .

Note that the unit balls have well-defined volume (i.e. 3-dimensional Lebesgue measure) and translation and rotations are very well behaved and preserve volume, while the subsets  $A_k$  are very pathological (i.e. non-Lebesgue-measurable).



Figure 31: Illustration of the Banach-Tarski paradox.<sup>68</sup>

⇒ Measure theory is the unifying ‘safe space’ for probability theory!

## A.2. Core Concepts

**Motivation A.2.1.** As discussed before in remark A.1.4, we want to define probability measures  $P$  on a space  $\mathcal{W}$ . We want them to follow (at least) these rules:

- i) *normalized*:  $P(\mathcal{W}) = 1$ ,  $P(\emptyset) = 0$ .
- ii) *complement*:  $P(A^c) = 1 - P(A)$  for  $A \subseteq \mathcal{W}$ .
- iii)  *$\sigma$ -additivity (aka countably additivity)*: For pairwise disjoint subsets  $A_n \subseteq \mathcal{W}$ ,  $n \in \mathbb{N}$ :

$$P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} P(A_n).$$

Such rules implicitly assume that  $P$  can measure the sets  $\mathcal{W}$  and  $\emptyset$ ; and that  $P$  can measure the complement  $A^c$  if it can measure  $A$ ; and that  $P$  can measure the (disjoint) union  $\bigcup_{n \in \mathbb{N}} A_n$  if it can measure each of the  $A_n$ .

As illustrated by the theorems A.1.5 and A.1.7, this is in general NOT possible to do for all subsets of  $\mathcal{W}$  (i.e. for all elements of the power set  $2^{\mathcal{W}}$ ).

This problem is solved and formalized by the notion of  $\sigma$ -algebras of subsets of the space  $\mathcal{W}$ .

**Definition A.2.2** ( $\sigma$ -algebras). Let  $\mathcal{W}$  be a set. A (non-empty) set  $\mathcal{B} \subseteq 2^{\mathcal{W}}$  of subsets  $A \subseteq \mathcal{W}$  is called a  $\sigma$ -algebra on  $\mathcal{W}$  if it satisfies the following rules:

- i) *empty set*:  $\emptyset \in \mathcal{B}$ ,
- ii) *complement*: If  $A \in \mathcal{B}$  then also:  $A^c := \mathcal{W} \setminus A \in \mathcal{B}$ ,
- iii) *countable union*: If  $A_n \in \mathcal{B}$  for all  $n \in \mathbb{N}$  then also:  $\bigcup_{n \in \mathbb{N}} A_n \in \mathcal{B}$ .

**Definition A.2.3** (Measurable spaces). A tuple  $(\mathcal{W}, \mathcal{B})$  of a set  $\mathcal{W}$  and a  $\sigma$ -algebra  $\mathcal{B}$  on  $\mathcal{W}$  is called measurable space.

<sup>68</sup>[https://en.wikipedia.org/wiki/Banach-Tarski\\_paradox](https://en.wikipedia.org/wiki/Banach-Tarski_paradox)

**Remark A.2.4** (Abuse of notation). *By abuse of notation we often just call  $\mathcal{W}$  a measurable space by implicitly assuming that it is endowed with a fixed  $\sigma$ -algebra, which we will indicate by  $\mathcal{B}_{\mathcal{W}}$  or  $\mathcal{B}(\mathcal{W})$  if needed. We will also just call a subsets  $A \subseteq \mathcal{W}$  measurable when we actually mean that  $A \in \mathcal{B}_{\mathcal{W}}$ .*

**Definition A.2.5** (Measures). *Let  $(\mathcal{W}, \mathcal{B})$  be a measurable space. A measure  $\mu$  on  $(\mathcal{W}, \mathcal{B})$  - by definition - is a mapping:*

$$\mu : \mathcal{B} \rightarrow \mathbb{R} \cup \{\infty\}, \quad D \mapsto \mu(D),$$

such that:

- i) non-negative:  $\forall A \in \mathcal{B}: \mu(A) \in [0, \infty]$ ,
- ii) empty set:  $\mu(\emptyset) = 0$ ,
- iii) countably additive (aka  $\sigma$ -additive): for all sequences  $A_n \in \mathcal{B}$ ,  $n \in \mathbb{N}$ , with  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ , we have:

$$\mu \left( \bigcup_{n \in \mathbb{N}} A_n \right) = \sum_{n \in \mathbb{N}} \mu(A_n).$$

**Definition A.2.6** (Probability/finite/ $\sigma$ -finite measures). *A measure  $\mu$  on  $(\mathcal{W}, \mathcal{B})$  is called:*

1. probability measure if  $\mu(\mathcal{W}) = 1$ .
2. finite measure if  $\mu(\mathcal{W}) < \infty$ .
3.  $\sigma$ -finite measure if there are  $D_n \in \mathcal{B}$ ,  $n \in \mathbb{N}$ , with  $\mu(D_n) < \infty$  and  $\mathcal{W} = \bigcup_{n \in \mathbb{N}} D_n$ .

**Definition A.2.7** (Measure spaces/probability spaces). *A triple  $(\mathcal{W}, \mathcal{B}, \mu)$  consisting of a measurable space  $(\mathcal{W}, \mathcal{B})$  and a measure  $\mu$  on  $(\mathcal{W}, \mathcal{B})$  is called measure space (and probability space if  $\mu$  is a probability measure).*

*Again, by abuse of notation, we often omit the  $\sigma$ -algebra in the notation and call  $(\mathcal{W}, \mu)$  a measure space, probability space, resp.*

**Definition A.2.8** (Measurable mappings). *Let  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  and  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}})$  be two measurable spaces and  $f : \mathcal{W} \rightarrow \mathcal{Z}$  be a mapping. We call  $f$  a  $\mathcal{B}_{\mathcal{W}}$ - $\mathcal{B}_{\mathcal{Z}}$ -measurable mapping (or just measurable for short) if for all  $B \in \mathcal{B}_{\mathcal{Z}}$  the pre-image  $f^{-1}(B)$  is an element of  $\mathcal{B}_{\mathcal{W}}$ . In formulas:*

$$\forall B \in \mathcal{B}_{\mathcal{Z}} : f^{-1}(B) \in \mathcal{B}_{\mathcal{W}}.$$

*Remember the definition of pre-image:  $f^{-1}(B) := \{w \in \mathcal{W} \mid f(w) \in B\}$ .*

**Definition A.2.9** (Push-forward measure). Let  $X : (\mathcal{W}, \mathcal{B}_{\mathcal{W}}) \rightarrow (\mathcal{X}, \mathcal{B}_{\mathcal{X}})$  be measurable and  $\mu$  a measure on  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$ . Then we define the push-forward measure (aka image measure) of  $\mu$  via:

$$(X_*\mu)(A) := \mu^X(A) := \mu_X(A) := \mu(X)(A) := \mu(X \in A) := \mu(X^{-1}(A))$$

for all  $A \in \mathcal{B}_{\mathcal{X}}$ . If  $\mu$  is a probability distribution then the push-forward measure  $\mu(X)$  is also called the (distributional) law of  $X$ .

**Definition A.2.10** (Random variables). A measurable mapping:

$$X : (\mathcal{W}, \mathcal{B}_{\mathcal{W}}, P) \rightarrow (\mathcal{X}, \mathcal{B}_{\mathcal{X}})$$

that starts from a probability space is also called random variable.

The main point is that the map  $X$  comes with its own distribution  $P^X$ . We often just say: “Let  $X$  be a random variable with distribution  $P^X = \dots$ ”, where  $P^X$  is then specified, e.g. to be a Gaussian or a categorical distribution, etc.

**Definition A.2.11** (Null sets). Let  $(\mathcal{W}, \mathcal{B}, \mu)$  be a measure space. A subset  $M \subseteq \mathcal{W}$  is called  $\mu$ -null or  $\mu$ -zero set if there exists a set  $N \in \mathcal{B}$  with  $M \subseteq N$  and  $\mu(N) = 0$ .

**Definition A.2.12** (Almost surely/almost all). Let  $(\mathcal{X}, \mathcal{B}, \mu)$  be a measure space and  $f, g : \mathcal{X} \rightarrow \mathcal{Z}$  a measurable map. We write  $f =_{\mu} g$  or say  $f = g$   $\mu$ -almost-surely (a.s.) or  $f(x) = g(x)$  for  $\mu$ -almost-all  $x \in \mathcal{X}$  if:

$$\{x \in \mathcal{X} \mid f(x) \neq g(x)\} \quad \text{is a } \mu\text{-null set.}$$

Similarly, for  $f \leq_{\mu} g$ , etc..

More generally, we say that a condition  $C$  about points  $x \in \mathcal{X}$  holds  $\mu$ -almost-surely or for  $\mu$ -almost-all  $x \in \mathcal{X}$  if the set of points where the condition does not hold is  $\mu$ -null, i.e.:

$$\{x \in \mathcal{X} \mid \neg C(x)\} \quad \text{is a } \mu\text{-null set.}$$

### A.3. Default Choices for Sigma-Algebras

In this subsection we want to highlight what kind of default  $\sigma$ -algebras we will assume on different types of spaces and on spaces constructed from others.

**Remark A.3.1** (Discrete spaces). If  $\mathcal{W}$  is countable (i.e. either finite or countably infinite, e.g. like  $\mathbb{Z}$ ,  $\mathbb{Q}$  or  $\mathbb{N}$  or  $\{1, \dots, N\}$ ) then we will always implicitly assume that  $\mathcal{W}$  is endowed with the power set  $\sigma$ -algebra:  $\mathcal{B}_{\mathcal{W}} = 2^{\mathcal{W}}$  (unless stated otherwise).

**Definition A.3.2** ( $\sigma$ -algebra generated by a set of subsets). Let  $\mathcal{W}$  be a set and  $\mathcal{A} \subseteq 2^{\mathcal{W}}$  be any non-empty set of subsets of  $\mathcal{W}$ . Then we can define the  $\sigma$ -algebra generated by  $\mathcal{A}$ :

$$\sigma(\mathcal{A}) := \bigcap_{\substack{\mathcal{B} \subseteq 2^{\mathcal{W}} \\ \mathcal{A} \subseteq \mathcal{B} \\ \mathcal{B} \text{ } \sigma\text{-algebra on } \mathcal{W}}} \mathcal{B},$$

as the intersection of all  $\sigma$ -algebras  $\mathcal{B}$  on  $\mathcal{W}$  that contain  $\mathcal{A}$ . Note that the set  $\sigma(\mathcal{A})$  really is a well-defined  $\sigma$ -algebra on  $\mathcal{W}$ .  $\sigma(\mathcal{A})$  is thus - by definition - the smallest  $\sigma$ -algebra on  $\mathcal{W}$  that contains  $\mathcal{A}$ .

**Definition A.3.3** (Borel  $\sigma$ -algebra on topological spaces). Let  $(\mathcal{W}, \mathcal{O})$  be a topological space with set of open subsets  $\mathcal{O}$  then the Borel  $\sigma$ -algebra of  $(\mathcal{W}, \mathcal{O})$  is defined as the smallest  $\sigma$ -algebra that contains all open (and thus also all closed) subsets:

$$\mathcal{B}_{(\mathcal{W}, \mathcal{O})} := \sigma(\mathcal{O}).$$

We will always implicitly assume that every topological space is endowed with its Borel  $\sigma$ -algebra (unless stated otherwise).

**Remark A.3.4.** *Caution:* Other choices of  $\sigma$ -algebras for topological spaces used in the literature are the Baire  $\sigma$ -algebra, which is generated by the zero sets of all continuous functions, or the  $\sigma$ -algebra generated only by its closed (countably) compact sets, or the  $\sigma$ -algebra of all (Radon-)universally measurable subsets.

**Lemma A.3.5** (Borel  $\sigma$ -algebra on  $\mathbb{R}^D$ ). The Borel  $\sigma$ -algebra of  $\mathbb{R}^D$  is generated by the cubes:

$$\mathcal{B}_{\mathbb{R}^D} = \sigma(\{[a_1, b_1] \times \cdots \times [a_D, b_D] \mid a_d, b_d \in \mathbb{Q}, a_d \leq b_d, d = 1, \dots, D\}).$$

**Definition/Lemma A.3.6** ( $\sigma$ -algebras induced by mappings). Let  $f : \mathcal{W} \rightarrow \mathcal{Z}$  be any mapping.

1. Let  $\mathcal{B}_{\mathcal{Z}}$  be a  $\sigma$ -algebra on  $\mathcal{Z}$ . Then the pull-back  $\sigma$ -algebra defined via:

$$f^* \mathcal{B}_{\mathcal{Z}} := \{f^{-1}(C) \mid C \in \mathcal{B}_{\mathcal{Z}}\}$$

is the smallest  $\sigma$ -algebra  $\mathcal{B}_{\mathcal{W}}$  that makes  $f$   $\mathcal{B}_{\mathcal{W}}$ - $\mathcal{B}_{\mathcal{Z}}$ -measurable.

2. Let  $\mathcal{B}_{\mathcal{W}}$  be a  $\sigma$ -algebra on  $\mathcal{W}$ . Then the push-forward  $\sigma$ -algebra defined via:

$$f_* \mathcal{B}_{\mathcal{W}} := \{C \subseteq \mathcal{Z} \mid f^{-1}(C) \in \mathcal{B}_{\mathcal{W}}\}$$

is the biggest  $\sigma$ -algebra  $\mathcal{B}_{\mathcal{Z}}$  that makes  $f$   $\mathcal{B}_{\mathcal{W}}$ - $\mathcal{B}_{\mathcal{Z}}$ -measurable.

**Definition A.3.7** (Product  $\sigma$ -algebra). Let  $(\mathcal{X}_i, \mathcal{B}_i)$  be measurable spaces,  $i \in I$ . Then the product space  $\prod_{i \in I} \mathcal{X}_i$  is endowed with the smallest  $\sigma$ -algebra such that for every  $j \in I$  the projection map:

$$\text{pr}_j : \prod_{i \in I} \mathcal{X}_i \rightarrow \mathcal{X}_j, \quad (x_i)_{i \in I} \mapsto x_j,$$

is measurable. We use the symbols  $\bigotimes_{i \in I} \mathcal{B}_i$  for this product  $\sigma$ -algebra. In symbols:

$$\bigotimes_{i \in I} \mathcal{B}_i := \sigma \left( \bigcup_{i \in I} \text{pr}_i^* \mathcal{B}_i \right).$$

We will always implicitly assume that every product space is endowed with this product  $\sigma$ -algebra (unless stated otherwise).

**Definition A.3.8** (Subspace  $\sigma$ -algebra). *Let  $(\mathcal{W}, \mathcal{B})$  be a measurable space and  $\mathcal{Z} \subseteq \mathcal{W}$  be a subset. Then the subspace  $\sigma$ -algebra  $\mathcal{B}|_{\mathcal{Z}}$  on  $\mathcal{Z}$  is the smallest  $\sigma$ -algebra that makes the inclusion map  $\mathcal{Z} \rightarrow \mathcal{W}$  measurable. More concretely:*

$$\mathcal{B}|_{\mathcal{Z}} := \{B \cap \mathcal{Z} \mid B \in \mathcal{B}\}.$$

*We will always assume that subsets are endowed with the subspace  $\sigma$ -algebra (unless it is ambiguous or stated otherwise).*

**Definition A.3.9** (Disjoint union  $\sigma$ -algebra). *Let  $(\mathcal{X}_i, \mathcal{B}_i)$  be measurable spaces,  $i \in I$ , considered to be pairwise disjoint. Then the disjoint union  $\sigma$ -algebra on the disjoint union  $\coprod_{i \in I} \mathcal{X}_i$  is the biggest  $\sigma$ -algebra  $\mathcal{B}_{\sqcup}$  such that all inclusion maps  $\mathcal{X}_i \rightarrow \coprod_{i \in I} \mathcal{X}_i$  are measurable. In symbols:*

$$\mathcal{B}_{\sqcup} := \left\{ E \subseteq \coprod_{i \in I} \mathcal{X}_i \mid \forall i \in I : E \cap \mathcal{X}_i \in \mathcal{B}_i \right\}.$$

**Definition A.3.10** ( $\sigma$ -algebra on the space of all probability measures). *Let  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  be a measurable space. We denote the space of all probability measures on  $(\mathcal{W}, \mathcal{B}_{\mathcal{W}})$  by:*

$$\mathcal{P}(\mathcal{W}) := \{P \mid P \text{ is probability measure on } (\mathcal{W}, \mathcal{B}_{\mathcal{W}})\}.$$

*We endow  $\mathcal{P}(\mathcal{W})$  with the smallest  $\sigma$ -algebra  $\mathcal{B}_{\mathcal{P}(\mathcal{W})}$  such that all evaluation maps:*

$$\text{ev}_D : \mathcal{P}(\mathcal{W}) \rightarrow [0, 1], \quad P \mapsto P(D)$$

*are measurable for  $D \in \mathcal{B}_{\mathcal{W}}$ . In symbols:*

$$\mathcal{B}_{\mathcal{P}(\mathcal{W})} := \sigma \left( \bigcup_{D \in \mathcal{B}_{\mathcal{W}}} \text{ev}_D^* \mathcal{B}_{[0,1]} \right).$$

*We will always assume that the space of probability measures  $\mathcal{P}(\mathcal{W})$  is endowed with this  $\sigma$ -algebra (unless stated otherwise).*

## A.4. Standard Measurable Spaces

**Definition A.4.1** (Standard measurable space). *A measurable space  $(\mathcal{W}, \mathcal{B})$  is called standard measurable space (aka standard Borel space) if it is measurably isomorphic to either:*

1. *a finite measurable space  $\{1, \dots, M\}$  for some  $M \in \mathbb{N}$  endowed with the power set  $\sigma$ -algebra  $2^{\{1, \dots, M\}}$ , or:*
2. *the countably infinite space  $\mathbb{N}$  endowed with the power set  $\sigma$ -algebra  $2^{\mathbb{N}}$ , or:*
3. *the unit interval  $[0, 1]$  endowed with its Borel  $\sigma$ -algebra:*

$$\mathcal{B}_{[0,1]} = \sigma(\{[a, b] \mid a, b \in [0, 1] \cap \mathbb{Q}, a \leq b\}).$$

'Measurably isomorphic' means that there is a measurable mapping that has a measurable inverse.

**Theorem A.4.2** (Kuratowski et al.). 1. Every Borel subset of any complete metric space that has a countable dense subset is a standard measurable space in its Borel  $\sigma$ -algebra (e.g.  $\mathbb{Q}^D$  is countable and dense in  $\mathbb{R}^D$ ).

2. Two standard measurable spaces  $\mathcal{X}$  and  $\mathcal{Y}$  are measurably isomorphic iff their cardinalities  $|\mathcal{X}|, |\mathcal{Y}|$  are equal (e.g.  $\mathbb{R}^D \cong [0, 1]$ ).

3. Countable disjoint unions and countable direct products of standard measurable spaces are standard measurable spaces.

4. If  $\mathcal{W}$  is standard measurable space then the space of its probability measures  $\mathcal{P}(\mathcal{W})$  is also a standard measurable space.

**Example A.4.3.** Examples of standard measurable spaces are:  $\mathbb{R}^D, \mathbb{Q}, \mathbb{Z}, \mathbb{N}, \{1, \dots, M\}, [0, 1]$ , topological manifolds, countable CW-complexes, every Borel set of any separable complete metric space.

## A.5. Measure Integrals

The construction of the measure integral  $\int f d\mu$  follows in several steps.

**Construction A.5.1** (Measure integral). Let  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}}, \mu)$  be a measure space.

1. Indicator functions: For  $A \in \mathcal{B}_{\mathcal{X}}$  put:

$$\int \mathbb{1}_A d\mu := \mu(A).$$

2. Simple functions: For a simple function  $g : \mathcal{X} \rightarrow \mathbb{R}$  given by:

$$g(x) = \sum_{n=1}^N a_n \cdot \mathbb{1}_{A_n}(x),$$

where  $A_n \in \mathcal{B}_{\mathcal{X}}$  and  $a_n \in \mathbb{R}, n = 1, \dots, N$ , we define:

$$\int g d\mu := \sum_{n=1}^N a_n \cdot \mu(A_n).$$

3. Non-negative measurable functions: Let  $h : \mathcal{X} \rightarrow [0, \infty]$  be a non-negative measurable function then we define:

$$\int h d\mu := \sup_{0 \leq g \leq h} \int g d\mu \quad \in [0, \infty],$$

where the supremum is running over all non-negative simple functions  $g$  that are smaller or equal to  $h$ .

4. *Measurable functions with well-defined integral: Let  $f : \mathcal{X} \rightarrow \bar{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$  be a measurable function. We then can write  $f = f_+ - f_-$  with:*

$$f_+ := \max(f, 0) \geq 0, \quad f_- := \max(-f, 0) \geq 0.$$

*If at least one of  $\int f_+ d\mu$ ,  $\int f_- d\mu$  is finite (i.e.  $< \infty$ ) we can then define:*

$$\int f d\mu := \int f_+ d\mu - \int f_- d\mu \quad \in [-\infty, \infty].$$

*The only case where we cannot properly define the integral is for measurable functions  $f : \mathcal{X} \rightarrow \bar{\mathbb{R}}$  where both integrals:  $\int f_+ d\mu = \infty$  and  $\int f_- d\mu = \infty$  are infinite, because of the “ $\infty - \infty = ?$ ” problem.*

**Remark A.5.2** (Riemann integral vs. measure integral). *The construction of the Riemann integral (RI) and the measure integral (MI) differ only in a few points:*

1. *RI uses (infinitesimal) interval length on  $x$ -axis, while MI uses the measure content (which is also the interval length in case of the Lebesgue measure).*
2. *RI decomposes the  $x$ -axis, while MI decomposes the  $y$ -axis.*
3. *RI uses limits for the integration boundaries to integrate to infinity (if convergent), while MI takes difference of integrals (to infinity) of  $f_+$  and  $f_-$  (if difference well-defined).*
4. *RI integrates in direction  $a$  to  $b$ , while MI integrates interval  $[a, b]$  in an undirected fashion.*
5. *If a function is Riemann integrable (RI) (e.g. continuous) on interval  $[a, b]$  then it is also Lebesgue integrable (MI) with the same integral value.*

**Definition A.5.3** (Integrable functions). *Let  $(\mathcal{X}, \mathcal{B}, \mu)$  be a measure space. A measurable function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is called  $\mu$ -integrable if:*

$$\int |f| d\mu < \infty.$$

**Theorem A.5.4** (Properties of the integral). *Let  $(\mathcal{X}, \mathcal{B}, \mu)$  be a measure space and  $f, g : \mathcal{X} \rightarrow \mathbb{R}$   $\mu$ -integrable measurable functions.*

1. *For  $A \in \mathcal{B}$  we have:  $\int \mathbb{1}_A d\mu = \mu(A)$ .*
2. *Linearity: If  $a, b \in \mathbb{R}$  then  $a \cdot f + b \cdot g$  is also  $\mu$ -integrable and we have:*

$$\int (a \cdot f + b \cdot g) d\mu = a \cdot \int f d\mu + b \cdot \int g d\mu.$$



3. *Triangle inequality:*

$$\left| \int f d\mu \right| \leq \int |f| d\mu < \infty.$$

4. *If  $f \geq_{\mu} 0$  then:  $\int f d\mu \geq 0$ , with equality iff  $f =_{\mu} 0$ .*

5. *Monotonicity: If  $f \geq_{\mu} g$  then:  $\int f d\mu \geq \int g d\mu$ , with equality iff  $f =_{\mu} g$ .*

6. *If  $\int f d\mu < \infty$  then  $f <_{\mu} \infty$ .*

7. *The measure integral satisfies monotone convergence, dominated convergence, Fubini theorems, etc. (see literature).*

*Note, we use  $=_{\mu}$  and  $\geq_{\mu}$  to indicate that this property is (only) allowed to fail on a  $\mu$ -null set.*

**Definition A.5.5** (Expectation value). *Let  $(\mathcal{W}, \mathcal{B}, P)$  be a probability space and  $X : \mathcal{W} \rightarrow \mathbb{R}$  be a measurable function with well-defined integral. Then its expectation value (w.r.t.  $P$ ) is defined to be:*

$$\mathbb{E}[X] := \int X dP.$$

**Example A.5.6.** *Let  $\mathcal{X}$  be a measurable space and  $f : \mathcal{X} \rightarrow \mathbb{R}$  a measurable function and  $\mathcal{W} \subseteq \mathcal{X}$  a countable subset.*

1. *Dirac measure. Let  $w \in \mathcal{X}$  be a point. We define the Dirac measure  $\delta_w$  centered at  $w$  via:*

$$\delta_w(A) := \mathbb{1}_A(w),$$

*for all measurable  $A \subseteq \mathcal{X}$ . Furthermore, we have:*

$$\mathbb{E}[f] = \int f(x) \delta_w(dx) = f(w).$$

*This holds because:  $f(x) = f(w)$  for  $\delta_w$ -almost-all  $x \in \mathcal{X}$ . Let's prove the right equality more formally:*

*Proof.* Consider:

$$B := f^{-1}(f(w)) = \{x \in \mathcal{X} \mid f(x) = f(w)\} \ni w.$$

Since  $f$  is measurable and  $\{f(w)\} \in \mathcal{B}_{\mathbb{R}}$  we also have  $B \in \mathcal{B}_{\mathcal{X}}$ . We then have the decomposition:

$$\begin{aligned} f(x) &= f(x) \cdot \mathbb{1}_B(x) + f(x) \cdot \mathbb{1}_{B^c}(x) \\ &= f(w) \cdot \mathbb{1}_B(x) + f(x) \cdot \mathbb{1}_{B^c}(x). \end{aligned}$$

Since  $w \notin B^c$  we get  $\delta_w(B^c) = 0$  and thus  $\int f(x) \cdot \mathbb{1}_{B^c}(x) \delta_w(dx) = 0$ . Together we get:

$$\begin{aligned} \int f(x) \delta_w(dx) &= \int (f(w) \cdot \mathbb{1}_B(x) + f(x) \cdot \mathbb{1}_{B^c}(x)) \delta_w(dx) \\ &= f(w) \cdot \int \mathbb{1}_B(x) \delta_w(dx) + \underbrace{\int f(x) \cdot \mathbb{1}_{B^c}(x) \delta_w(dx)}_{=0} \\ &= f(w) \cdot \delta_w(B) \\ &= f(w). \end{aligned}$$

□

2. Discrete distributions. Consider a discrete probability distribution  $P$  supported on the countable subset  $\mathcal{W} \subseteq \mathcal{X}$ . Let  $p$  be its mass function. We then can write the corresponding probability measure  $P$  on  $\mathcal{X}$  as:

$$P = \sum_{w \in \mathcal{W}} p(w) \cdot \delta_w.$$

For measurable  $A \subseteq \mathcal{X}$  we then have:

$$P(A) = \sum_{w \in \mathcal{W}} p(w) \cdot \delta_w(A) = \sum_{w \in \mathcal{W} \cap A} p(w).$$

Furthermore, we get:

$$\mathbb{E}[f] = \int f(x) P(dx) = \int f(x) \sum_{w \in \mathcal{W}} p(w) \cdot \delta_w(dx) = \sum_{w \in \mathcal{W}} f(w) \cdot p(w).$$

## A.6. Densities/Derivatives

**Definition A.6.1.** Let  $(\mathcal{X}, \mathcal{B})$  be a measure space and  $\mu, \nu$  two measures on it. We say that  $\nu$  has a density w.r.t.  $\mu$  if there exists a non-negative measurable function  $f : \mathcal{X} \rightarrow [0, \infty]$  such that for all  $A \in \mathcal{B}$ :

$$\nu(A) = \int \mathbb{1}_A \cdot f d\mu =: \int_A f d\mu.$$

Such a density does not always exist. If a density exists then it is essentially unique, in the sense that two such densities would only differ on a  $\mu$ -null set. We often use the notation:  $f = \frac{d\nu}{d\mu}$  and call it ‘the’ density or (Radon-Nikodým) derivative of  $\nu$  w.r.t.  $\mu$ .

**Proposition A.6.2.** Let  $(\mathcal{X}, \mathcal{B})$  be a measure space and  $\mu, \nu, \kappa$  three measures on it and  $g : \mathcal{X} \rightarrow \bar{\mathbb{R}}$  is either a  $\nu$ -integrable or non-negative measurable function.

1. If  $\nu$  has a density w.r.t.  $\mu$  then we have:

$$\int g d\nu = \int g \cdot \frac{d\nu}{d\mu} d\mu.$$

2. (Conic) linearity: If  $\kappa$  has a density w.r.t.  $\mu$  and  $\nu$  has a density w.r.t.  $\mu$  and  $a, b \geq 0$  then  $a \cdot \kappa + b \cdot \nu$  has a density w.r.t.  $\mu$  and we have:

$$\frac{d(a \cdot \kappa + b \cdot \nu)}{d\mu}(x) = a \cdot \frac{d\kappa}{d\mu}(x) + b \cdot \frac{d\nu}{d\mu}(x)$$

for  $\mu$ -almost-all  $x \in \mathcal{X}$ .

3. Chain rule: If  $\nu$  has a density w.r.t.  $\mu$  and  $\mu$  has a density w.r.t.  $\kappa$  then also  $\nu$  has a density w.r.t.  $\kappa$  and we have:

$$\frac{d\nu}{d\kappa}(x) = \frac{d\nu}{d\mu}(x) \cdot \frac{d\mu}{d\kappa}(x)$$

for  $\kappa$ -almost-all  $x \in \mathcal{X}$ .

4. Inverse: If  $\nu$  has a density w.r.t.  $\mu$  and  $\mu$  has a density w.r.t.  $\nu$  then we have:

$$\frac{d\nu}{d\mu}(x) = \left( \frac{d\mu}{d\nu}(x) \right)^{-1}$$

for  $\mu$ -almost-all  $x \in \mathcal{X}$ . We can make in this context the (somewhat arbitrary) choice to put:  $0^{-1} := \infty$ .

**Definition A.6.3** (Absolute continuity). Let  $\mu, \nu$  be two measures on a measurable space  $(\mathcal{X}, \mathcal{B})$ . We say that  $\nu$  is absolute continuous w.r.t.  $\mu$ , in symbols:

$$\nu \ll \mu,$$

if for every  $A \in \mathcal{B}$  with  $\mu(A) = 0$  also  $\nu(A) = 0$  holds, in short, if:

$$\mu(A) = 0 \quad \implies \quad \nu(A) = 0.$$

**Theorem A.6.4** (Radon-Nikodým, see [Kle20] Cor. 7.34). Let  $(\mathcal{X}, \mathcal{B}, \mu)$  be a  $\sigma$ -finite measure space and  $\nu$  another measure on  $(\mathcal{X}, \mathcal{B})$ . Then the following two statements are equivalent:

1.  $\nu$  has a density w.r.t.  $\mu$ .
2.  $\nu$  is absolute continuous w.r.t.  $\mu$ .

**Theorem A.6.5** (Besicovitch density theorem, [Fre15] 472D). Let  $\mu$  be a Radon measure on  $\mathbb{R}^D$  (e.g. any finite or probability measure or the Lebesgue measure, see A.8.1) and  $f : \mathbb{R}^D \rightarrow \bar{\mathbb{R}}$  be any (locally)  $\mu$ -integrable function. Then we have for  $\mu$ -almost-all  $x \in \mathbb{R}^D$ :

1.  $\lim_{\varepsilon \rightarrow 0} \frac{1}{\mu(B_\varepsilon(x))} \int_{B_\varepsilon(x)} f(z) \mu(dz) = f(x).$
2.  $\lim_{\varepsilon \rightarrow 0} \frac{1}{\mu(B_\varepsilon(x))} \int_{B_\varepsilon(x)} |f(z) - f(x)| \mu(dz) = 0.$

Here  $B_\varepsilon(x)$  denote the closed balls of radius  $\varepsilon > 0$  centered at  $x$  (in Euclidean norm). The above, in particular, holds for the density  $f = \frac{d\nu}{d\mu}$  of another measure  $\nu$  w.r.t.  $\mu$ :

$$\lim_{\varepsilon \rightarrow 0} \frac{\nu(B_\varepsilon(x))}{\mu(B_\varepsilon(x))} = \frac{d\nu}{d\mu}(x),$$

for  $\mu$ -almost-all  $x \in \mathbb{R}^D$ .

## A.7. Conditional Expectation

You may be familiar with the conditional expectation for discrete random variables  $X, Y$ :

$$\begin{aligned} \mathbb{E}[X|Y = y] &= \sum_{x \in \mathcal{X}} x \cdot P(X = x|Y = y) &&= \sum_{x \in \mathcal{X}} x \cdot \frac{P(X = x, Y = y)}{P(Y = y)} \\ &= \frac{\sum_{x \in \mathcal{X}} x \cdot P(X = x, Y = y)}{P(Y = y)}, \end{aligned}$$

and for real-valued random variables  $X, Y$  with positive and continuous joint density  $p(x, y)$ :

$$\mathbb{E}[X|Y = y] = \int_{\mathcal{X}} x \cdot p(x|Y = y) dx = \int_{\mathcal{X}} x \cdot \frac{p(x, y)}{p(y)} dx = \frac{\int_{\mathcal{X}} x \cdot p(x, y) dx}{p(y)}.$$

The following construction generalizes this notion:

**Definition A.7.1** (Conditional expectation). *Let  $(\mathcal{W}, P)$  be a probability space and  $X : \mathcal{W} \rightarrow \mathbb{R}$ ,  $Y : \mathcal{W} \rightarrow \mathcal{Y}$  be two random variables with either  $\mathbb{E}[|X|] < \infty$  or  $X \geq 0$  a.s.*

1. *The conditional expectation of  $X$  given  $Y = y$  is defined via:*

$$\mathbb{E}[X|Y = y] := \mathbb{E}[X_+|Y = y] - \mathbb{E}[X_-|Y = y] \quad \in \bar{\mathbb{R}},$$

where  $X_\pm := \max(\pm X, 0) \geq 0$  and:

$$\mathbb{E}[X_\pm|Y = y] := \frac{dE_\pm}{dP^Y}(y),$$

is the Radon-Nikodym derivative/density w.r.t.  $P^Y$  of the following measure on  $\mathcal{Y}$ :

$$E_\pm(B) := \mathbb{E}[X_\pm \cdot \mathbb{1}_B(Y)] = \int x \cdot \mathbb{1}_B(y) dP^{(X_\pm, Y)}(x, y).$$

One can easily see that  $E_\pm \ll P^Y$  and that the densities exist by the Radon-Nikodym theorem.

2. The conditional expectation of  $X$  given  $Y$  is then the measurable map defined via:

$$\mathbb{E}[X|Y] : \mathcal{W} \rightarrow \bar{\mathbb{R}}, \quad w \mapsto \mathbb{E}[X|Y](w) := \mathbb{E}[X|Y = Y(w)] = \mathbb{E}[X|Y = y]|_{y=Y(w)},$$

i.e. the composition of  $Y$  with the measurable map  $y \mapsto \mathbb{E}[X|Y = y]$ .

**Remark A.7.2.** The construction from above also works with a measure  $\mu$  such that  $\mu^Y$  is  $\sigma$ -finite (instead of  $P$ ) since we only need to guarantee the existence of the Radon-Nikodym derivative.

**Notation A.7.3.** Let  $\mathcal{W}, \mathcal{Z}, \mathcal{Y}$  be measurable spaces and  $Z : \mathcal{W} \rightarrow \mathcal{Z}$  and  $Y : \mathcal{W} \rightarrow \mathcal{Y}$  be measurable maps. We write:

$$Z \lesssim Y$$

if there exists a measurable function  $F : \mathcal{Y} \rightarrow \mathcal{Z}$  such that  $Z = F \circ Y$ ; in other words if  $Z$  is a deterministic (measurable) function of  $Y$ , i.e.:  $Z = F(Y)$ .

If  $\mu$  is a measure on  $\mathcal{W}$  we also write:

$$Z \lesssim_{\mu} Y$$

if there exists a measurable map  $F$  such that  $Z = F(Y)$   $\mu$ -almost-surely.

**Theorem A.7.4.** Let  $(\mathcal{W}, P)$  be a probability space and  $X, T : \mathcal{W} \rightarrow \mathbb{R}, Y : \mathcal{W} \rightarrow \mathcal{Y}, Z : \mathcal{W} \rightarrow \mathcal{Z}$  be random variables with  $\mathbb{E}[|X|] < \infty$  (or as long as we do not run into the “ $\infty - \infty = ?$ ” problem). Then we have the following properties:

1.  $\mathbb{E}[X|Y]$  is the unique real valued random variable  $Z$  (up to  $P$ -null set) such that:
  - a)  $Z \lesssim_P Y$  and:
  - b) for all measurable  $B \subseteq \mathcal{Y}$ :

$$\mathbb{E}[Z \cdot \mathbb{1}_B(Y)] = \mathbb{E}[X \cdot \mathbb{1}_B(Y)].$$

2. For all real valued random variables  $Z \lesssim_P Y$  with  $\mathbb{E}[|Z \cdot X|] < \infty$  we have:

$$\mathbb{E}[Z \cdot X|Y] = Z \cdot \mathbb{E}[X|Y] \quad P\text{-a.s.}$$

3. Linearity: For all  $a, b \in \mathbb{R}$  we have:

$$\mathbb{E}[a \cdot X + b \cdot T|Y] = a \cdot \mathbb{E}[X|Y] + b \cdot \mathbb{E}[T|Y] \quad P\text{-a.s.}$$

4. Constants:  $\mathbb{E}[1|Y] = 1$   $P$ -a.s.

5. Constant maps: If  $Y$  is a constant map then:  $\mathbb{E}[X|Y] = \mathbb{E}[X]$   $P$ -a.s.

6. Independence (see 2.5.23): If  $X \perp\!\!\!\perp Y$  then:  $\mathbb{E}[X|Y] = \mathbb{E}[X]$   $P$ -a.s.

7. Deterministic dependence: If  $X \lesssim_P Y$  then:  $\mathbb{E}[X|Y] = X$   $P$ -a.s.

8. *Monotonicity: If  $X \geq T$   $P$ -a.s. then we have:*

$$\mathbb{E}[X|Y] \geq \mathbb{E}[T|Y] \quad P\text{-a.s.}$$

9. *Jensen inequality: Let  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  be convex then we have:*

$$\varphi(\mathbb{E}[X|Y]) \leq \mathbb{E}[\varphi(X)|Y] \quad P\text{-a.s.}$$

10. *Triangle inequality:  $|\mathbb{E}[X|Y]| \leq \mathbb{E}[|X||Y]$   $P$ -a.s.*

11. *Tower rule: If  $Y \preceq Z$  then:*

$$\mathbb{E}[\mathbb{E}[X|Y]|Z] = \mathbb{E}[\mathbb{E}[X|Z]|Y] = \mathbb{E}[X|Y] \quad P\text{-a.s.}$$

12. *Tower rule, special case:*

$$\mathbb{E}[\mathbb{E}[X|Y]|Y, Z] = \mathbb{E}[\mathbb{E}[X|Y, Z]|Y] = \mathbb{E}[X|Y] \quad P\text{-a.s.}$$

13. *Monotone convergence, dominated convergence, etc. (see literature).*

## A.8. The Lebesgue Measure

**Definition A.8.1** (The Lebesgue (outer) measure). *The Lebesgue (outer) measure  $\lambda^D$  on  $\mathbb{R}^D$  is given for subsets  $A \subseteq \mathbb{R}^D$  via:*

$$\lambda^D(A) := \inf \left\{ \sum_{n \in \mathbb{N}} \text{vol}^D([a^{(n)}, b^{(n)}]) \mid A \subseteq \bigcup_{n \in \mathbb{N}} [a^{(n)}, b^{(n)}] \right\},$$

where the infimum is running over sequences of  $D$ -dimensional cubes:

$$[a^{(n)}, b^{(n)}] = [a_1^{(n)}, b_1^{(n)}] \times \cdots \times [a_D^{(n)}, b_D^{(n)}],$$

with  $a^{(n)} = (a_1^{(n)}, \dots, a_D^{(n)})$ ,  $b^{(n)} = (b_1^{(n)}, \dots, b_D^{(n)}) \in \mathbb{R}^D$ ,  $a_d^{(n)} \leq b_d^{(n)}$  for  $d = 1, \dots, D$ ,  $n \in \mathbb{N}$ , that jointly cover  $A$ , where the  $D$ -dimensional volume is given by:

$$\text{vol}^D([a^{(n)}, b^{(n)}]) := (b_1^{(n)} - a_1^{(n)}) \cdots (b_D^{(n)} - a_D^{(n)}), \quad \text{vol}^D(\emptyset) := 0.$$

**Theorem A.8.2** (The Lebesgue measure). *The Lebesgue measure  $\lambda^D$ , when restricted to the Borel- $\sigma$ -algebra of  $\mathbb{R}^D$ , is the unique measure on  $\mathbb{R}^D$  that satisfies:*

$$\lambda^D([a, b]) = \text{vol}^D([a, b]),$$

for all  $D$ -dimensional cubes  $[a, b]$ . If the dimension is clear from the context we might just write  $\lambda$  for  $\lambda^D$ .

**Theorem A.8.3.** Let  $\lambda$  be the Lebesgue measure on the interval  $[a, b]$ . Let  $f : [a, b] \rightarrow \mathbb{R}$  be a Riemann integrable function (e.g. a continuous function) then  $f$  is also  $\lambda$ -integrable and we have:

$$\int_a^b f(x) dx = \int_{[a,b]} f(x) \lambda(dx).$$

**Theorem A.8.4** (Fundamental theorem of calculus). Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a measurable function such that  $\int_{[a,b]} |f| d\lambda < \infty$  for  $a, b \in \mathbb{R}$ . For fixed  $c \in \mathbb{R}$  define  $F : [c, \infty) \rightarrow \mathbb{R}$  via:

$$F(x) := \int_{(c,x]} f d\lambda.$$

Then  $F$  is differentiable in  $\lambda$ -almost-all  $x \in \mathbb{R}$  and for those points we have:

$$F'(x) = f(x).$$

## A.9. Transformation Rules

**Theorem A.9.1** (General integral transformation). Let  $(\mathcal{W}, \mu)$  be a measure space and  $X : \mathcal{W} \rightarrow \mathcal{X}$  and  $F : \mathcal{X} \rightarrow \mathbb{R}$  be measurable. Then we have:

$$\int F(X) d\mu = \int F d(X_*\mu),$$

if either side is well-defined. Written in longer form this is:

$$\int F(X(w)) \mu(dw) = \int F(x) (X_*\mu)(dx).$$

**Theorem A.9.2** (Push-forward of densities). Let  $(\mathcal{W}, \mu)$  be a measure space and  $\nu$  another measure on  $\mathcal{W}$ . Let  $\varphi : \mathcal{W} \rightarrow \mathcal{Y}$  be a measurable mapping such that  $\varphi_*\mu$  is  $\sigma$ -finite. If  $\nu$  has a density w.r.t.  $\mu$  then the push-forward measure  $\varphi_*\nu$  has a density w.r.t.  $\varphi_*\mu$  given as follows:

$$\frac{d(\varphi_*\nu)}{d(\varphi_*\mu)}(y) = \mathbb{E}_\mu \left[ \frac{d\nu}{d\mu} \middle| \varphi = y \right] = \int \frac{d\nu}{d\mu}(w) \mu(dw | \varphi = y),$$

for  $\varphi_*\mu$ -almost-all  $y \in \mathcal{Y}$ , where the conditional integral  $\mathbb{E}_\mu$  is constructed the same way as the conditional expectation but using the  $\sigma$ -finite measure  $\varphi_*\mu$ .

If, furthermore,  $\varphi$  is a measurable isomorphism then we get:

$$\frac{d(\varphi_*\nu)}{d(\varphi_*\mu)}(y) = \frac{d\nu}{d\mu}(\varphi^{-1}(y))$$

for  $\varphi_*\mu$ -almost-all  $y \in \mathcal{Y}$ .

**Theorem A.9.3** (Transformation formula for the Lebesgues measure). Let  $\varphi : \mathbb{R}^D \rightarrow \mathbb{R}^D$  be a continously differentiable bijection of  $\mathbb{R}^D$  (or of open/closed subsets therein)

with Jacobian  $\varphi'(x)$  at point  $x$ . Let  $\lambda$  be the Lebesgue measure on  $\mathbb{R}^D$ . Then  $\varphi_*\lambda$  is absolute continuous w.r.t.  $\lambda$  with density given by:

$$\frac{d(\varphi_*\lambda)}{d\lambda}(y) = |\det \varphi'(\varphi^{-1}(y))|^{-1}$$

for all  $y \in \mathbb{R}^D$  (or in that open/closed subset, and  $= 0$  outside).

**Corollary A.9.4** (Transformation of (probability) densities w.r.t. the Lebesgue measure). Let the setting be like in A.9.3. Let  $\nu$  be a (probability) measure on  $\mathbb{R}^D$  with (probability) density  $p$  w.r.t.  $\lambda$ . Then  $\varphi_*\nu$  also has a (probability) density w.r.t.  $\lambda$ , which is then given by:

$$\frac{d(\varphi_*\nu)}{d\lambda}(y) = \frac{d(\varphi_*\nu)}{d(\varphi_*\lambda)}(y) \cdot \frac{d(\varphi_*\lambda)}{d\lambda}(y) = p(\varphi^{-1}(y)) \cdot |\det \varphi'(\varphi^{-1}(y))|^{-1}.$$

**Theorem A.9.5** (A bit more general, [Fre15] Cor. 263F, 262F(b)). Let  $\mathcal{X} \subseteq \mathbb{R}^D$  be a measurable set and  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^D$  an injective Lipschitz function. Let  $\mathcal{X}' \subseteq \mathcal{X}$  be the set of points  $x$  at which  $\varphi$  has a derivative  $\varphi'(x)$  relative to  $\mathcal{X}$ <sup>69</sup>. Then we have:

1.  $\mathcal{X} \setminus \mathcal{X}'$  is a  $\lambda$ -null set.
2.  $|\det \varphi'| : \mathcal{X}' \rightarrow [0, \infty)$  is measurable.
3.  $\varphi(\mathcal{X}) \subseteq \mathbb{R}^D$  is a measurable set.
4.  $\lambda(\varphi(\mathcal{X})) = \int_{\mathcal{X}} |\det \varphi'(x)| d\lambda(x)$ .
5. For every real-valued function  $g$  defined on a subset  $\mathcal{Y} \subseteq \varphi(\mathcal{X})$  we have:

$$\int_{\varphi(\mathcal{X})} g(y) d\lambda(y) = \int_{\mathcal{X}} g(\varphi(x)) \cdot |\det \varphi'(x)| d\lambda(x),$$

if either integral is defined in  $[-\infty, \infty]$  and provided we interpret  $g(\varphi(x)) \cdot |\det \varphi'(x)| := 0$  if  $\varphi(x) \notin \mathcal{Y}$  and  $|\det \varphi'(x)| = 0$ .

**Remark A.9.6** (Transformation rule for discrete measures). Let  $\mathcal{X}$  be a measurable space and  $\mu$  be a discrete (probability) measure on  $\mathcal{X}$  supported on the countable discrete subset  $\mathcal{W} \subseteq \mathcal{X}$  with mass function given by:

$$m(x) = \frac{d\mu}{d\#\mathcal{W}}(x),$$

where  $\#\mathcal{W}$  is the counting measure w.r.t.  $\mathcal{W}$  given by:  $\#\mathcal{W}(A) := \#(\mathcal{W} \cap A)$ . Let  $\varphi : \mathcal{X} \rightarrow \mathcal{Y}$  be a measurable map. Then  $\varphi_*\mu$  is a discrete measure supported on  $\varphi(\mathcal{W})$  with mass function/density:

$$\frac{d\varphi_*\mu}{d\#\varphi(\mathcal{W})}(y) = \sum_{w \in \varphi^{-1}(y) \cap \mathcal{W}} m(w).$$

<sup>69</sup>We say that  $\varphi$  is differentiable relative to  $\mathcal{X}$  at  $x \in \mathcal{X}$  if there exists  $\varphi'(x) \in \mathbb{R}^{D \times D}$  such that for every  $\epsilon > 0$  there exists a  $\delta > 0$  such that for all  $y \in \mathcal{X}$  with  $\|y - x\| < \delta$  we have that:  $\|\varphi(y) - \varphi(x) - \varphi'(x) \cdot (y - x)\| \leq \epsilon \cdot \|y - x\|$ . Note that in this definition such a derivative  $\varphi'(x)$  does not need to be unique.



**Example A.9.7** (Linear transformation of Gaussian distributions).

**Example A.9.8** (Density of Chi-square distributions).

## A.10. Measure Extension Theorems

**Theorem A.10.1** (Measure extension theorem, see [Kle20] Thm. 1.53, Thm. 1.36). *Let  $\mathcal{A}$  be a ring (e.g. an algebra) of subsets of a set  $\mathcal{X}$ . Let  $\mu : \mathcal{A} \rightarrow [0, \infty)$  be a (finitely) additive set function with  $\mu(\emptyset) = 0$  that is also  $\emptyset$ -continuous:*

$$\inf_{n \in \mathbb{N}} \mu(A_n) = 0, \quad (71)$$

for all non-increasing sequences  $(A_n)_{n \in \mathbb{N}}$  with  $A_n \in \mathcal{A}$  and  $A_{n+1} \subseteq A_n$  for all  $n \in \mathbb{N}$  and  $\bigcap_{n \in \mathbb{N}} A_n = \emptyset$ .

Then there exists a unique  $\sigma$ -finite measure  $\nu : \sigma(\mathcal{A}) \rightarrow [0, \infty]$  such that  $\nu(A) = \mu(A)$  for all  $A \in \mathcal{A}$ .

**Theorem A.10.2** (Ionescu-Tulcea extension theorem, see [IT49, Lam87]). *Let  $I$  be an arbitrary set (not necessarily countable) and  $(\mathcal{X}_i, \mathcal{B}_i)$ ,  $i \in I$ , measurable spaces. For subsets  $J \subseteq I$  we put:*

$$\mathcal{X}_J := \prod_{j \in J} \mathcal{X}_j, \quad \mathcal{B}_J := \bigotimes_{j \in J} \mathcal{B}_j, \quad (72)$$

the product space endowed with its product  $\sigma$ -algebra. Now assume that we have a probability measure  $\mu_J$  on  $(\mathcal{X}_J, \mathcal{B}_J)$  for every finite subset  $J \subseteq I$  such that:

1. for every finite subsets  $L \subseteq J \subseteq I$  we have:  $\text{pr}_{L,*} \mu_J = \mu_L$ ,
2. for every finite subset  $J \subseteq I$  and  $i \in I \setminus J$  there exists a Markov kernel:  $\mu_{i|J} : \mathcal{X}_J \dashrightarrow \mathcal{X}_i$  such that:

$$\mu_{\{i\} \dot{\cup} J} = \mu_{i|J} \otimes \mu_J. \quad (73)$$

Then there exists a probability measure  $\mu_I$  on  $(\mathcal{X}_I, \mathcal{B}_I)$  such that for every finite subset  $J \subseteq I$  we have:

$$\text{pr}_{J,*} \mu_I = \mu_J. \quad (74)$$

*Proof.* We first put, with  $\text{pr}_J : \mathcal{X}_I \rightarrow \mathcal{X}_J$  the canonical projections:

$$\mathcal{A} := \bigcup_{\substack{J \subseteq I \\ \#J < \infty}} \text{pr}_J^* \mathcal{B}_J = \{ \text{pr}_J^{-1}(B) \subseteq \mathcal{X}_I \mid J \subseteq I, \#J < \infty, B \in \mathcal{B}_J \}. \quad (75)$$

Then, per definition,  $\mathcal{B}_I = \sigma(\mathcal{A})$ . Furthermore,  $\mathcal{A}$  is an algebra of subsets of  $\mathcal{X}_I$ . Indeed, let  $A_1, A_2 \in \mathcal{A}$  then  $A_l \in \text{pr}_{J_l}^* \mathcal{B}_{J_l}$  for some finite subsets  $J_l \subseteq I$ ,  $l = 1, 2$ . Then  $J := J_1 \cup J_2$

is also a finite subset of  $I$  and we have:  $A_1, A_2 \in \text{pr}_J^* \mathcal{B}_J$ . So,  $A_l = \text{pr}_J^{-1}(B_l)$  for some  $B_l \in \mathcal{B}_J$ ,  $l = 1, 2$ . This then shows that:

$$A_l^c = \text{pr}_J^{-1}(B_l^c) \in \text{pr}_J^* \mathcal{B}_J \subseteq \mathcal{A}, \quad (76)$$

$$A_1 \cup A_2 = \text{pr}_J^{-1}(B_1 \cup B_2) \in \text{pr}_J^* \mathcal{B}_J \subseteq \mathcal{A}, \quad (77)$$

$$A_1 \cap A_2 = \text{pr}_J^{-1}(B_1 \cap B_2) \in \text{pr}_J^* \mathcal{B}_J \subseteq \mathcal{A}. \quad (78)$$

It is also clear that:  $\mathcal{X}_I, \emptyset \in \mathcal{A}$ . So,  $\mathcal{A}$  is an algebra of subsets of  $\mathcal{X}_I$ .

We can now define the set function  $\mu : \mathcal{A} \rightarrow [0, 1]$  via:

$$\mu(A) := \mu_J(B), \quad \text{for } A = \text{pr}_J^{-1}(B), \quad B \in \mathcal{B}_J. \quad (79)$$

This is well-defined because of the condition:  $\text{pr}_{L,*} \mu_J = \mu_L$  for finite subsets  $L \subseteq J \subseteq I$ .

It is also clear that  $\mu$  is additive. Indeed, if  $A_1, A_2 \in \mathcal{A}$  are disjoint then  $A_l = \text{pr}_J^{-1}(B_l)$  for some finite subset  $J \subseteq I$  and some disjoint  $B_l \in \mathcal{B}_J$ ,  $l = 1, 2$ . The additivity of  $\mu_J$  then shows the additivity of  $\mu$ :

$$\mu(A_1 \dot{\cup} A_2) = \mu_J(B_1 \dot{\cup} B_2) = \mu_J(B_1) + \mu_J(B_2) = \mu(A_1) + \mu(A_2). \quad (80)$$

To apply the extension theorem [A.10.1](#) it is left to check that  $\mu$  is  $\emptyset$ -continuous on  $\mathcal{A}$ . For this, and, by way of contradiction, consider a non-increasing sequence  $A_n \in \mathcal{A}$ ,  $n \in \mathbb{N}$ , with  $\bigcap_{n \in \mathbb{N}} A_n = \emptyset$  and  $\inf_{n \in \mathbb{N}} \mu(A_n) > \epsilon > 0$ . We can assume that  $A_n = \text{pr}_{J_n}^{-1}(B_n)$  with  $B_n \in \mathcal{B}_{J_n}$  with the inclusion of finite subsets:  $J_n \subseteq J_{n+1} \subseteq I$  for all  $n \in \mathbb{N}$ . We totally order the countable set  $\bigcup_{n \in \mathbb{N}} J_n$  such that  $k < l$  if  $k \in J_n$  and  $l \in J_{n+1} \setminus J_n$ .

We introduce the following abbreviations for  $n \in \mathbb{N}$ :

$$\mathcal{Y}_n := \mathcal{X}_{J_n \setminus J_{n-1}}, \quad \mathcal{Y}_{\leq n} := \prod_{l=1}^n \mathcal{Y}_l, \quad \mathcal{Y}_c := \mathcal{X}_{I \setminus \bigcup_{n \in \mathbb{N}} J_n}, \quad (81)$$

$$\mu_{n|<n} := \bigotimes_{k \in J_n \setminus J_{n-1}} \mu_{k|\{l \in J_n \mid l < k\} \cup J_{n-1}}, \quad \mu_{\leq n} := \mu_{J_n}. \quad (82)$$

Then  $\mathcal{Y} := \mathcal{Y}_c \times \prod_{n \in \mathbb{N}} \mathcal{Y}_n = \mathcal{X}_I$ . We also put:

$$h_n(y) := g_n(y_{\leq n}) := \mathbb{1}_{B_n}(y_{\leq n}) = \mathbb{1}_{A_n}(y), \quad h(y) := \inf_{n \in \mathbb{N}} h_n(y). \quad (83)$$

By assumption,  $\bigcap_{n \in \mathbb{N}} A_n = \emptyset$ , we have that  $h(y) = 0$  for all  $y \in \mathcal{Y}$ . Since  $A_n \subseteq A_{n-1}$  we have for all  $n \in \mathbb{N}$ :

$$0 = h(y) \leq h_n(y) \leq h_{n-1}(y) \leq 1. \quad (84)$$

We define for  $k, n \in \mathbb{N}$ :

$$f_n^{(k)}(y_{\leq k}) := \int g_n(y_{k+1:n}, y_{\leq k}) \mu_{k+1:n|\leq k}(dy_{k+1:n} | y_{\leq k}), \quad (85)$$

$$f^{(k)}(y_{\leq k}) := \inf_{n \in \mathbb{N}} f_n^{(k)}(y_{\leq k}). \quad (86)$$

Note that we then also have for all  $k, n \in \mathbb{N}$  and  $y \in \mathbb{N}$ :

$$0 \leq f^{(k)}(y_{\leq k}) \leq f_n^{(k)}(y_{\leq k}) \leq f_{n-1}^{(k)}(y_{\leq k}) \leq 1. \quad (87)$$

We also put for  $k \in \mathbb{N}$ :

$$C_{\leq k} := \{y_{\leq k} \in \mathcal{Y}_{\leq k} \mid f^{(k)}(y_{\leq k}) > \epsilon\}. \quad (88)$$

By the above assumption we have for every  $n \in \mathbb{N}$ :

$$0 < \epsilon < \inf_{n \in \mathbb{N}} \mu(A_n) \quad (89)$$

$$= \inf_{n \in \mathbb{N}} \mu_{\leq n}(B_n) \quad (90)$$

$$= \inf_{n \in \mathbb{N}} \mathbb{E}_{\mu_{\leq n}}[g_n] \quad (91)$$

$$= \inf_{n \in \mathbb{N}} \int \int g_n(y_{2:n}, y_1) \mu_{2:n|1}(dy_{2:n}|y_1) \mu_1(dy_1) \quad (92)$$

$$= \inf_{n \in \mathbb{N}} \int f_n^{(1)}(y_1) \mu_1(dy_1) \quad (93)$$

$$= \int \inf_{n \in \mathbb{N}} f_n^{(1)}(y_1) \mu_1(dy_1) \quad (94)$$

$$= \int f^{(1)}(y_1) \mu_1(dy_1). \quad (95)$$

Where the integral and infimum can be interchanged because  $f^{(1)}$  is  $\mathcal{B}_1$ -measurable and the monotone convergence theorem, applied to  $1 - f_n^{(1)}$ . We see that:  $\mu_1(C_1) > 0$ . Otherwise,  $\int f^{(1)}(y_1) \mu_1(dy_1) \leq \epsilon$ , which would contradict the above sequence of inequalities.

Now inductively for  $k \in \mathbb{N}$  and  $y_{\leq k} \in C_{\leq k}$  we have:

$$0 < \epsilon < f^{(k)}(y_{\leq k}) \quad (96)$$

$$= \inf_{n \in \mathbb{N}} f_n^{(k)}(y_{\leq k}) \quad (97)$$

$$= \inf_{n \in \mathbb{N}} \int f_n^{(k+1)}(y_{k+1}, y_{\leq k}) \mu_{k+1|\leq k}(dy_{k+1}|y_{\leq k}) \quad (98)$$

$$= \int \inf_{n \in \mathbb{N}} f_n^{(k+1)}(y_{k+1}, y_{\leq k}) \mu_{k+1|\leq k}(dy_{k+1}|y_{\leq k}) \quad (99)$$

$$= \int f^{(k+1)}(y_{k+1}, y_{\leq k}) \mu_{k+1|\leq k}(dy_{k+1}|y_{\leq k}). \quad (100)$$

This shows that:  $\mu_{k+1|\leq k}(C_{\leq k+1}^{y_{\leq k}}|y_{\leq k}) > 0$  for  $y_{\leq k} \in C_{\leq k}$ . This means that we can inductively construct a  $y \in \mathcal{Y}$  with components:  $y_1 \in C_1$  and  $y_{k+1} \in C_{\leq k+1}^{y_{\leq k}}$  for  $k \in \mathbb{N}$ , and an arbitrary  $y_c \in \mathcal{Y}_c$ . This  $y$  then satisfies  $h_n(y) > \epsilon > 0$  for all  $n \in \mathbb{N}$  and thus  $h(y) = \inf_{n \in \mathbb{N}} h_n(y) \geq \epsilon > 0$ , which lies in contradiction to  $h(y) = 0$  for all  $y \in \mathcal{Y}$ .

This shows that  $\mu$  is  $\emptyset$ -continuous. It follows by the extension theorem A.10.1 that  $\mu$  has a unique extension to a probability measure to  $\sigma(\mathcal{A}) = \mathcal{B}_I$ . This shows the claim.  $\square$

**Corollary A.10.3** (Ionescu-Tulcea extension theorem for Markov kernels). *Let  $I$  be an arbitrary set (not necessarily countable) and  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}})$  and  $(\mathcal{X}_i, \mathcal{B}_i)$ ,  $i \in I$ , measurable spaces. For subsets  $J \subseteq I$  we put:*

$$\mathcal{X}_J := \prod_{j \in J} \mathcal{X}_j, \quad \mathcal{B}_J := \bigotimes_{j \in J} \mathcal{B}_j, \quad (101)$$

*the product space endowed with its product  $\sigma$ -algebra. Now assume that for every finite subset  $J \subseteq I$  we are given a Markov kernel:*

$$K_J(X_J|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_J,$$

*such that:*

1. *for every finite subsets  $L \subseteq J \subseteq I$  we have:*

$$K_J(X_L|Z) = K_L(X_L|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_L,$$

2. *for every finite subset  $J \subseteq I$  and  $i \in I \setminus J$  there exists a Markov kernel:*

$$K_{i|J}(X_i|X_J, Z) : \mathcal{X}_J \times \mathcal{Z} \dashrightarrow \mathcal{X}_{\{i\} \cup J},$$

*such that:*

$$K_{i|J}(X_i|X_J, Z) \otimes K_J(X_J|Z) = K_{\{i\} \cup J}(X_i, X_J|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_{\{i\} \cup J}.$$

*Then there exists a Markov kernel:*

$$K(X_I|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_I,$$

*such that for every finite subset  $J \subseteq I$  we have:*

$$K(X_J|Z) = K_J(X_J|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_J.$$

*Proof.* For every  $z \in \mathcal{Z}$  we can apply the Ionescu-Tulcea extension theorem [A.10.2](#) separately and get a probability measure  $K(X_I|Z = z)$  on  $\mathcal{B}_I$  such that for every finite subset  $J \subseteq I$  we have:

$$K(X_J|Z = z) = K_J(X_J|Z = z).$$

We are left to check that the map:

$$K(X_I|Z) : \mathcal{Z} \rightarrow \mathcal{P}(\mathcal{X}_I), \quad z \mapsto K(X_I|Z = z),$$

is measurable. By Dynkin's lemma and the definition of the product  $\sigma$ -algebra  $\mathcal{B}_I$  this only needs to be checked on sets  $A_J \in \mathcal{B}_J$  for finite subsets  $J \subseteq I$ . Since we have:

$$K(X_I \in \text{pr}_J^{-1}(A)|Z = z) = K(X_J \in A|Z = z) = K_J(X_J \in A|Z = z),$$

and  $z \mapsto K_J(X_J \in A|Z = z)$  is measurable for finite  $J \subseteq I$ , the claim follows.  $\square$

**Corollary A.10.4** (Kolmogorov extension theorem for Markov kernels). *Let  $I$  be an arbitrary set (not necessarily countable)  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}})$  a measurable space and  $(\mathcal{X}_i, \mathcal{B}_i)$ ,  $i \in I$ , standard measurable spaces. For subsets  $J \subseteq I$  we put:*

$$\mathcal{X}_J := \prod_{j \in J} \mathcal{X}_j, \quad \mathcal{B}_J := \bigotimes_{j \in J} \mathcal{B}_j, \quad (102)$$

*the product space endowed with its product  $\sigma$ -algebra. Now assume that for every finite subset  $J \subseteq I$  we are given a Markov kernel:*

$$K_J(X_J|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_J,$$

*such that for every finite subsets  $L \subseteq J \subseteq I$  we have:*

$$K_J(X_L|Z) = K_L(X_L|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_L.$$

*Then there exists a Markov kernel:*

$$K(X_I|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_I,$$

*such that for every finite subset  $J \subseteq I$  we have:*

$$K(X_J|Z) = K_J(X_J|Z) : \mathcal{Z} \dashrightarrow \mathcal{X}_J.$$

*Proof.* This directly follows from Ionescu-Tulcea extension theorem for Markov kernels [A.10.3](#) and the fact that on standard measurable spaces we always have conditional Markov kernels by the disintegration theorem [2.4.16](#).  $\square$

## References

- [ARSZ05] R. Ayesha Ali, Thomas S. Richardson, Peter Spirtes, and Jiji Zhang. Towards characterizing Markov equivalence classes for directed acyclic graphs with latent variables. In *Proceedings of the 21th Conference on Uncertainty in Artificial Intelligence (UAI-05)*, pages 10–17, 2005.
- [BFPM21] Stephan Bongers, Patrick Forré, Jonas Peters, and Joris M. Mooij. Foundations of structural causal models with cycles and latent variables. *Annals of Statistics*, 49(5):2885–2915, 2021. doi:10.1214/21-AOS2064.
- [BJvdV17] Fetsje Bijma, Marianne Jonker, and Aad van der Vaart. *An Introduction to Mathematical Statistics*. Amsterdam University Press, 2017.
- [BM18] Stephan Bongers and Joris M. Mooij. From random differential equations to structural causal models: the stochastic case. *arXiv.org preprint*, arXiv:1803.08784v2 [cs.AI], March 2018. URL: <https://arxiv.org/abs/1803.08784v2>.
- [BP94a] A. Balke and J. Pearl. Counterfactual probabilities: computational methods, bounds and applications. In *Proceedings of the Tenth International Conference on Uncertainty in Artificial Intelligence (UAI-94)*, pages 46–54, 1994.
- [BP94b] A. Balke and J. Pearl. Probabilistic evaluation of counterfactual queries. In *Proceedings of the 12th Conference on Artificial Intelligence*, volume 1, pages 230–237. MIT Press, 1994.
- [BTP14] E. Bareinboim, J. Tian, and J. Pearl. Recovering from selection bias in causal and statistical inference. In *AAAI*, 2014.
- [BvDM21] Tineke Blom, Mirthe M. van Diepen, and Joris M. Mooij. Conditional independences and causal relations implied by sets of equations. *Journal of Machine Learning Research*, 22(178):1–62, 2021. URL: <http://jmlr.org/papers/v22/20-863.html>.
- [CB17] J. D. Correa and E. Bareinboim. Causal effect identification by adjustment under confounding and selection biases. In *AAAI*, 2017.
- [CD17] Panayiota Constantinou and A. Philip Dawid. Extended conditional independence and applications in causal inference. *The Annals of Statistics*, pages 2618–2653, 2017.
- [CH11] Tom Claassen and Tom Heskes. A logical characterization of constraint-based causal discovery. In *Proceedings of the 27th Annual Conference on Uncertainty in Artificial Intelligence (UAI 2011)*, pages 135–144, 2011.

- [CMH13] Tom Claassen, Joris M. Mooij, and Tom Heskes. Learning sparse causal models is not NP-hard. In Ann Nicholson and Padhraic Smyth, editors, *Proceedings of the 29th Annual Conference on Uncertainty in Artificial Intelligence (UAI-13)*, pages 172–181. AUAI Press, 2013. URL: <http://auai.org/uai2013/prints/papers/121.pdf>.
- [Coo97] Gregory F. Cooper. A simple constraint-based algorithm for efficiently mining observational databases for causal relationships. *Data Mining and Knowledge Discovery*, 1(2):203–224, 1997.
- [CTB18] J. D. Correa, J. Tian, and E. Bareinboim. Generalized adjustment under confounding and selection biases. In *AAAI*, 2018.
- [Dar53] George Darmois. Analyse générale des liaisons stochastiques: etude particulière de l’analyse factorielle linéaire. *Revue de l’Institut international de statistique*, pages 2–8, 1953.
- [Daw79] A. Philip Dawid. Conditional independence in statistical theory. *Journal of the Royal Statistical Society: Series B (Methodological)*, 41(1):1–15, 1979.
- [Daw80] A. Philip Dawid. Conditional independence for statistical operations. *The Annals of Statistics*, pages 598–617, 1980.
- [Daw01] A. Philip Dawid. Separoids: a mathematical framework for conditional independence and irrelevance. *Annals of Mathematics and Artificial Intelligence*, 32(1-4):335–372, 2001.
- [Daw02] A. Philip Dawid. Influence diagrams for causal modelling and inference. *International Statistical Review*, 70:161–189, 2002.
- [DM83] Claude Dellacherie and Paul-André Meyer. *Probability and Potential B. Theory of Martingales*. Number 72 in North-Holland Mathematics Studies. Elsevier Science, 1983. Translated from French by J. P. Wilson.
- [FM17] Patrick Forré and Joris M. Mooij. Markov properties for graphical models with cycles and latent variables. *arXiv.org preprint*, arXiv:1710.08775 [math.ST], 2017. URL: <https://arxiv.org/abs/1710.08775>.
- [FM18] Patrick Forré and Joris M. Mooij. Constraint-based causal discovery for non-linear structural causal models with cycles and latent confounders. In *Proceedings of the 34th Annual Conference on Uncertainty in Artificial Intelligence (UAI-2018)*, 2018.
- [FM20] Patrick Forré and Joris M. Mooij. Causal calculus in the presence of cycles, latent confounders and selection bias. In *Proceedings of the 35th Annual Conference on Uncertainty in Artificial Intelligence (UAI-2019)*, volume 115, pages 71–80. PMLR, 2020. URL: <http://proceedings.mlr.press/v115/forre20a.html>.

- [For21] Patrick Forré. Transitional conditional independence. *arXiv.org preprint*, arXiv:2104.11547 [math.ST], 2021. URL: <https://arxiv.org/abs/2104.11547>.
- [Fre15] David H. Fremlin. *Measure Theory*, volume 1-6. Torres Fremlin, 2000-2015. URL: <https://www1.essex.ac.uk/maths/people/fremlin/mt.htm>.
- [GP95] D. Galles and J. Pearl. Testing identifiability of causal effects. In *UAI*, 1995.
- [Haa43] T. Haavelmo. The statistical implications of a system of simultaneous equations. *Econometrica*, 11(1):1–12, 1 1943.
- [HV06] Y. Huang and M. Valtorta. Pearl’s calculus of intervention is complete. In *UAI*, 2006.
- [HV08] Y. Huang and M. Valtorta. On the completeness of an identifiability algorithm for semi-Markovian models. *Annals of Mathematics and Artificial Intelligence*, 54(4):363–408, 2008.
- [IT49] Cassius T. Ionescu Tulcea. Mesures dans les espaces produits. *Atti Accad. Naz. Lincei Rend*, 7:208–211, 1949.
- [Kal17] Olav Kallenberg. *Random Measures, Theory and Applications*, volume 77 of *Probability Theory and Stochastic Modelling*. Springer, 2017.
- [Kec95] Alexander S. Kechris. *Classical Descriptive Set Theory*, volume 156 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1995.
- [Kle20] Achim Klenke. *Probability Theory—A Comprehensive Course*. Universitext. Springer, 3rd edition, 2020.
- [Lam87] Charles W. Lamb. A Comparison of Methods for Constructing Probability Measures on Infinite Product Spaces. *Canadian Mathematical Bulletin*, 30(3):282–285, 1987.
- [LDLL90] S.L. Lauritzen, A.P. Dawid, B.N. Larsen, and H.-G. Leimer. Independence properties of directed Markov fields. *Networks*, 20(5):491–505, 1990.
- [Man06] Subramani Mani. *A Bayesian Local Causal Discovery Framework*. PhD thesis, University of Pittsburg, March 2006. URL: <http://d-scholarship.pitt.edu/10181/>.
- [MC20] Joris M. Mooij and Tom Claassen. Constraint-based causal discovery using partial ancestral graphs in the presence of cycles. In Jonas Peters and David Sontag, editors, *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI-20)*, volume 124, pages 1159–1168. PMLR, 8 2020. URL: <http://proceedings.mlr.press/v124/m-mooij20a/m-mooij20a-supp.pdf>.



- [Mes12] Franz H. Messerli. Chocolate consumption, cognitive function, and Nobel laureates. *N Engl J. Med.*, 367:1562–1564, 2012. doi:[10.1056/NEJMon1211064](https://doi.org/10.1056/NEJMon1211064).
- [MMC20] Joris M. Mooij, Sara Magliacane, and Tom Claassen. Joint causal inference from multiple contexts. *Journal of Machine Learning Research*, 21(99):1–108, 2020. URL: <http://jmlr.org/papers/v21/17-123.html>.
- [MN98] Charles F. Manski and Daniel S. Nagin. Bounding disagreements about treatment effects. *Sociological Methodology*, 28(1):99–137, 1998.
- [MP13] Charles F. Manski and John V. Pepper. Deterrence and the death penalty: Partial identification analysis using repeated cross sections. *Journal of Quantitative Criminology*, 29:123–141, 2013.
- [Par05] Kalyanapuram Rangachari Parthasarathy. *Probability Measures on Metric Spaces*, volume 352. American Mathematical Society, 2005.
- [PBM16] Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society, Series B*, 78(5):947–1012, 2016.
- [Pea93a] J. Pearl. Aspects of graphical models connected with causality. In *Proceedings of the 49th Session of the International Statistical Institute*, pages 391–401, 1993.
- [Pea93b] J. Pearl. Comment: Graphical models, causality, and intervention. *Statistical Science*, 8(3):266–269, 1993.
- [Pea95] J. Pearl. Causal diagrams for empirical research (with discussion). *Biometrika*, 82(4):669–710, 1995.
- [Pea09] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2nd edition, 2009.
- [PP85] Judea Pearl and Azaria Paz. *Graphoids: A Graph-based Logic for Reasoning about Relevance Relations*. University of California (Los Angeles). Computer Science Department, 1985.
- [PP10] J. Pearl and A. Paz. Confounding equivalence in causal inference. In *UAI*, 2010.
- [PTKM15] E. Perkovic, J. Textor, M. Kalisch, and M. Maathuis. A complete generalized adjustment criterion. In *UAI*, 2015.
- [RERS23] Thomas S. Richardson, Robin J. Evans, James M. Robins, and Ilya Shpitser. Nested Markov Properties for Acyclic Directed Mixed Graphs. *The Annals of Statistics*, 51(1):334–361, 2023.

- [Ric03] Thomas S. Richardson. Markov properties for acyclic directed mixed graphs. *Scandinavian Journal of Statistics*, 30(1):145–157, 2003.
- [RR13a] Thomas S. Richardson and James M. Robins. Single world intervention graphs: A primer. In *Second UAI Workshop on Causal Structure Learning, Bellevue, Washington*, 2013.
- [RR13b] Thomas S. Richardson and James M. Robins. Single world intervention graphs (SWIGs): A unification of the counterfactual and graphical approaches to causality. *Center for the Statistics and the Social Sciences, University of Washington Series. Working Paper*, 128(30), 2013.
- [RS02] Thomas S. Richardson and Peter Spirtes. Ancestral graph Markov models. *The Annals of Statistics*, 30(4):962–1030, August 2002.
- [SC99] Peter Spirtes and Gregory F. Cooper. An experiment in causal discovery using a pneumonia database. In David Heckerman and Joe Whittaker, editors, *Proceedings of the Seventh International Workshop on Artificial Intelligence and Statistics*, volume R2 of *Proceedings of Machine Learning Research*. PMLR, 1999.
- [SdWR10] I. Shpitser, T. J. Van der Weele, and J. M. Robins. On the validity of covariate adjustment for estimating causal effects. In *UAI*, 2010.
- [SGS00] Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, Prediction, and Search*. MIT Press, Cambridge, MA, 2nd edition, 2000.
- [SMR95] Peter Spirtes, Christopher Meek, and Thomas S. Richardson. Causal inference in the presence of latent variables and selection bias. In *Proceedings of the 11th Annual Conference on Uncertainty in Artificial Intelligence (UAI-95)*, pages 499–506. Morgan Kaufmann, 1995.
- [SMR99] Peter Spirtes, Christopher Meek, and Thomas Richardson. An algorithm for causal inference in the presence of latent variables and selection bias. In Clark Glymour and Gregory F. Cooper, editors, *Computation, Causation, and Discovery*, chapter 6, pages 211–252. MIT Press, Cambridge, MA, 1999.
- [SP06a] I. Shpitser and J. Pearl. Identification of conditional interventional distributions. In *UAI*, 2006.
- [SP06b] I. Shpitser and J. Pearl. Identification of joint interventional distributions in recursive semi-Markovian causal models. In *AAAI*, 2006.
- [Spi94] P. Spirtes. Conditional independence in directed cyclic graphical models for feedback. Technical Report CMU-PHIL-54, Carnegie Mellon University, 1994.

- [Spi95] P. Spirtes. Directed cyclic graphical representations of feedback models. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence (UAI-95)*, pages 499–506, 1995.
- [SPP<sup>+</sup>05] Karen Sachs, Omar Perez, Dana Pe’er, Douglas A. Lauffenburger, and Garry P. Nolan. Causal protein-signaling networks derived from multiparameter single-cell data. *Science*, 308, 2005.
- [SW60] Robert H. Strotz and H. O. A. Wold. Recursive vs. nonrecursive systems: An attempt at synthesis. *Econometrica*, 28(2):417–427, 1960.
- [Tia02] J. Tian. *Studies in Causal Reasoning and Learning*. PhD thesis, University of California, Los Angeles, 2002.
- [Tia04] J. Tian. Identifying conditional causal effects. In *UAI*, 2004.
- [TP02] J. Tian and J. Pearl. A general identification condition for causal effects. In *AAAI*, 2002.
- [Č82] Nikolai Nikolaevich Čencov. *Statistical Decision Rules and Optimal Inference*, volume 53 of *Translations of Mathematical Monographs*. American Mathematical Society, 1982. translated from Russian.
- [vdV98] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- [Ver93] Tom S. Verma. Graphical aspects of causal models. Technical Report R-191, Computer Science Department, University of California, Los Angeles, 1993.
- [vH48] Joanne Baptista van Helmont. *Ortus medicinae: Id est Initia physicae inaudita. Progressus medicinae novus, in morborum ultionem, ad vitam longam*. Apud Ludovicum Elzevirium, 1648.
- [VM19] Philip Versteeg and Joris M. Mooij. Boosting local causal discovery in high-dimensional expression data. In *2019 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 2599–2604, Los Alamitos, CA, USA, 11 2019. IEEE Computer Society. URL: <https://doi.ieeecomputersociety.org/10.1109/BIBM47256.2019.8983232>, doi:10.1109/BIBM47256.2019.8983232.
- [VP90] T. Verma and J. Pearl. Equivalence and synthesis of causal models. In L. Kanal M. Henrion, R. Shachter and J. Lemmer, editors, *Uncertainty in Artificial Intelligence (UAI-90)*, pages 220–227. Association for Uncertainty in AI, 1990.
- [Wag20] Stefan Wager. Stats 361: Causal inference. Technical report, Stanford University, 2020. URL: <https://web.stanford.edu/~swager/stats361.pdf>.

- [Wri21] S. Wright. Correlation and causation. *Journal of Agricultural Research*, 20:557–585, 1921.
- [Zha06] Jiji Zhang. *Causal Inference and Reasoning in Causally Insufficient Systems*. PhD thesis, Carnegie Mellon University, July 2006. URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.466.7206&rep=rep1&type=pdf>.
- [Zha08] Jiji Zhang. On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16-17):1873–1896, 2008.